



HAL
open science

Objective Bayesian analysis of Kriging models with anisotropic correlation kernel

Joseph Muré

► **To cite this version:**

Joseph Muré. Objective Bayesian analysis of Kriging models with anisotropic correlation kernel. Probability [math.PR]. Université Sorbonne Paris Cité, 2018. English. NNT : 2018USPCC069 . tel-02184403

HAL Id: tel-02184403

<https://theses.hal.science/tel-02184403v1>

Submitted on 16 Jul 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Université Sorbonne
Paris Cité



École Doctorale Paris Centre

THÈSE DE DOCTORAT

Discipline : Mathématiques

Thèse de l'Université Sorbonne Paris Cité
préparée à l'Université Paris Diderot

présentée par

Joseph MURÉ

Objective Bayesian analysis of Kriging models with anisotropic correlation kernel

dirigée par Josselin GARNIER

Soutenue publiquement à Paris le 5 octobre 2018 devant le jury composé de :

Pr Clémentine PRIEUR	Université Grenoble Alpes	présidente
Pr Kerrie MENGERSEN	Queensland University of Technology	rapporteuse
Pr Jean-Michel MARIN	Université de Montpellier	rapporteur
Pr Josselin GARNIER	École Polytechnique	directeur
Pr Stéphane BOUCHERON	Université Paris Diderot	examinateur
Dr Loïc LE GRATIET	EDF R&D Chatou	encadrant
Mme Anne DUTFOY	EDF R&D Saclay	encadrante

Thèse réalisée dans le cadre d'un Contrat de Formation par la Recherche à EDF.

Remerciements

Je remercie en premier lieu Josselin Garnier, mon directeur de thèse. Tout au long de ces trois années de travail, il a veillé à la rigueur de mes travaux, en étant pleinement disponible pour suivre leurs progrès, semaine après semaine. Il m'a laissé conduire mes recherches dans les directions qui m'intéressaient, ce dont je lui suis profondément reconnaissant.

Je remercie aussi Loïc Le Gratiet, mon premier encadrant à EDF, qui m'a fait confiance par deux fois, d'abord au cours du stage préliminaire que j'ai réalisé à EDF, puis pour effectuer ce travail de thèse. Je lui suis reconnaissant pour ses précieux conseils et ses nombreux encouragements.

Je remercie de même Anne Dutfoy qui, ayant rejoint l'encadrement de ma thèse en cours de route, a immédiatement saisi ses enjeux, m'a apporté son aide et m'a témoigné une grande bienveillance.

Je suis honoré que Kerrie Mengersen et Jean-Michel Marin soient rapporteurs de ma thèse. Leurs rapports très détaillés m'ont permis d'améliorer la clarté de ce travail, et leurs encouragements me poussent à le faire connaître. Je les remercie pour cela. *I am honored that Kerrie Mengersen and Jean-Michel Marin reviewed this dissertation. Their detailed recommendations allowed me to make this work clearer, and I will keep their encouragement in mind when disseminating it. I thank them for that.* Je remercie aussi Clémentine Prieur et Stéphane Boucheron, qui ont accepté d'être membres de mon jury de thèse.

Je veux remercier tous mes collègues et ex-collègues du département PRISME d'EDF R&D, pour l'ambiance de travail chaleureuse qu'ils font régner et qui est si essentielle pour mener à bien avec sérénité un travail de longue haleine comme une thèse de doctorat. Je tiens à remercier particulièrement Merlin Keller et Nicolas Bousquet, qui m'ont souvent donné de judicieux conseils en me faisant profiter de leur expérience.

Je remercie le Laboratoire de Probabilités, Statistique et Modélisation de Sorbonne Université pour son accueil et la mise à disposition de ses ressources.

Je veux remercier aussi les doctorants des départements PRISME et PERICLES d'EDF R&D, dont l'entraide est vraiment un atout précieux. Je pense tout particulièrement à ceux qui sont déjà partis, Nazih Benoumechiara, Xavier Yau, Tuan Dinh Trong, Thomas Browne, et notamment à nos vacances inoubliables à Aussois. Et à ceux qui sont toujours là, Thomas Galtier, Jérôme Stenger, Thomas Bittar, Pablo Pereira-Alvarez, Mouna Rifi, Léna Masson, Artémis Drakos, Aurélie Daanen, Déborah Fontaine.

Enfin, je tiens à remercier mes parents pour leur soutien constant, ainsi que Jérémie Genty, Arthur Soulié, Guillaume Metzler, amis fidèles sans qui je ne serais pas arrivé jusqu'à ce point.

Abstract

A recurring problem in surrogate modelling is the scarcity of available data which hinders efforts to estimate model parameters. The Bayesian paradigm offers an elegant way to circumvent the problem by describing knowledge of the parameters by a posterior probability distribution instead of a pointwise estimate. However, it involves defining a prior distribution on the parameter. In the absence of expert opinion, finding an adequate prior can be a trying exercise. The Objective Bayesian school proposes default priors for such situations, like the Berger-Bernardo reference prior. Such a prior was derived by Berger et al. [2001] for the Kriging surrogate model with isotropic covariance kernel. Directly extending it to anisotropic kernels poses theoretical as well as practical problems because the reference prior framework requires ordering the parameters. Any ordering would in this case be arbitrary. Instead, we propose an Objective Bayesian solution for Kriging models with anisotropic covariance kernels based on conditional reference posterior distributions. This solution is made possible by a theory of compromise between incompatible conditional distributions. The work is then shown to be compatible with Trans-Gaussian Kriging. It is applied to an industrial case with nonstationary data in order to derive Probability Of defect Detection (POD) by non-destructive tests in steam generator tubes of nuclear power plants.

Keywords: Incompatibility, Conditional distribution, Markov kernel, Optimal compromise, Kriging, Reference prior.

Résumé

Les métamodèles statistiques sont régulièrement confrontés au manque de données qui engendre des difficultés à estimer les paramètres. Le paradigme bayésien fournit un moyen élégant de contourner le problème en décrivant la connaissance que nous avons des paramètres par une loi de probabilité *a posteriori* au lieu de la résumer par une estimation ponctuelle. Cependant, ce paradigme nécessite de définir une loi *a priori* adéquate, ce qui est un exercice difficile en l'absence de jugement d'expert. L'école bayésienne objective propose des priors par défaut dans ce genre de situation tels que le prior de référence de Berger-Bernardo. Un tel prior a été calculé par Berger et al. [2001] pour le modèle de krigeage avec noyau de covariance isotrope. Une extension directe au cas des noyaux anisotropes poserait des problèmes théoriques aussi bien que pratiques car la théorie de Berger-Bernardo ne peut s'appliquer qu'à un jeu de paramètres ordonnés. Or dans ce cas de figure, tout ordre serait nécessairement arbitraire. Nous y substituons une solution bayésienne objective fondée sur les posteriors de référence conditionnels. Cette solution est rendue possible par une théorie du compromis entre lois conditionnelles incompatibles. Nous montrons en outre qu'elle est compatible avec le krigeage trans-gaussien. Elle est appliquée à un cas industriel avec des données non-stationnaires afin de calculer des Probabilités de Détection de défauts (POD de l'anglais *Probability Of Detection*) par tests non-destructifs dans les tubes de générateur de vapeur de centrales nucléaires.

Mots-clés : Incompatibilité, Loi conditionnelle, Noyau markovien, Compromis optimal, Krigeage, Prior de référence.

Contents

Contents	5
Introduction	7
I Tools	11
1 Kriging Overview	13
1.1 Introduction	13
1.2 Gaussian random processes	14
1.3 Mean square continuity and differentiability of Gaussian processes	20
1.4 Spectral representation	25
1.5 Examples of correlation kernels	31
1.6 Current Kriging-related research	33
2 Reference Prior Theory	35
2.1 Introduction	35
2.2 Basic idea	36
2.3 Full definition of the reference prior	41
2.4 Regular continuous case	49
2.5 Properties of reference priors	55
2.6 Examples	56
2.7 Reference priors for multiparametric models	59
3 Propriety of the reference posterior distribution in Gaussian Process regression	65
3.1 Introduction	65
3.2 Setting	66
3.3 Smoothness of the correlation kernel	68
3.4 Propriety of the reference posterior distribution	70
3.5 Conclusion	72
Appendices	74
3.A Algebraic facts	74
3.B Maclaurin series	77
3.C Spectral decomposition	81
3.D Asymptotic study of the correlation matrix Σ_θ	85
3.E Details of the proof of Theorem 3.9	88

II Compromise	93
4 Optimal compromise between incompatible conditional probability distributions	95
4.1 Introduction	95
4.2 Optimal compromise: a general theory	96
4.3 Discussion of the notion of compromise	102
4.4 Conclusion	108
III Application	109
5 Application of the Optimal Compromise to Simple Kriging models with Matérn correlation kernels	111
5.1 Introduction	111
5.2 Optimal compromise between Objective Posterior conditional distributions in Gaussian Process regression	113
5.3 Comparisons between the MLE and MAP estimators	117
5.4 Comparison of the predictive distributions associated with the estimators (MLE and MAP) and the full posterior distribution	120
5.5 Conclusion and Perspectives	124
Appendices	126
5.A Proofs of Section 5.2	126
6 A Comprehensive Bayesian Treatment of the Universal Kriging model with Matérn correlation kernels	139
6.1 Introduction	140
6.2 Analytical treatment of the location-scale parameters	141
6.3 Reference prior on a one-dimensional θ	145
6.4 The Gibbs reference posterior on a multi-dimensional θ	145
6.5 Comparison of the predictive performance of the full-Bayesian approach versus MLE and MAP plug-in approaches	148
6.6 Conclusion	156
Appendices	159
6.A Matérn kernels	159
6.B Proofs of the existence of the Gibbs reference posterior	159
7 Trans-Gaussian Kriging in a Bayesian framework: a case study	167
7.1 Introduction	167
7.2 Probability Of Detection (POD)	169
7.3 An Objective Bayesian outlook to Trans-Gaussian Kriging	170
7.4 Industrial Application	177
7.5 Conclusion	181
Conclusion	183
Bibliography	185

Introduction

Objective

This thesis was funded by EDF and the ANR in order to solve a specific problem of emulation of a computer experiment in a situation with drastic budgetary constraints.

Emulation of computer codes has become in the last decades an object of considerable attention from the statistical community [Santner et al., 2003]. Surrogate models are statistical models designed to represent uncertainty about the result of experiments when the computational budget is limited. Linear regression, neural networks, polynomial chaos and Kriging are examples of surrogate models.

The problem EDF was facing was computing Probability Of defect Detection (POD) curves for a non-destructive test based on eddy-currents. A computer code was available to simulate such a procedure under various parameter choices. Because of the cost involved in making this code run, it was desirable to design a surrogate model that could be used in its stead. Because of the intrinsic Uncertainty Quantification capability of the Kriging model, which can easily provide prediction intervals or – critically for this application – probabilities of reaching a threshold, it was a natural choice. An apparent difficulty, however, was that the data collected seemed to defy the traditional Kriging assumption that a stationary Gaussian process could adequately represent them. There seemed to be a need to define non-Gaussian, non-stationary Kriging procedures.

Contributions of the thesis

It turned out, however, that a simple transformation of the data could tackle these problems, so there was no need to try any of the sophisticated procedures available in the literature. The true issue, which was not so easily solved, was that the relative lack of data (100 points dispersed in a 9-dimensional space) made estimating the Kriging hyperparameters difficult. This might not have been significant if they had not had a tremendous impact on prediction.

In the absence of a convincing estimator, it seemed reasonable to try to embed hyperparameter uncertainty into the prediction. And in this regard, the Bayesian approach seemed adequate. This framework treats the parameter as a random variable rather than as an unknown quantity. The associated probability distribution does not express actual randomness but is a quantification of an opinion about the value of the parameter.

Although philosophically coherent, this subjectivist view could not be used in this particular instance because there was no prior opinion that could be encoded by a Bayesian prior dis-

tribution. For the sake of convenience, it was necessary to find a generic, default approach in the absence of relevant information. The Berger-Bernardo [Bernardo, 1979a, Berger and Bernardo, 1992, Bernardo, 2005, Berger et al., 2009] “reference prior” theory provided a quasi-systematic approach that was quite appealing. It relies on a formally defined criterion of noninformativity aiming to “let the data speak for themselves”.

The reference prior was derived by Berger et al. [2001] for the Kriging model with isotropic correlation kernels – which require only one correlation parameter. They provided a proof, which we show is valid for a subset of all considered correlation kernels. Moreover, we contribute a proof for the complementary subset.

Aside from this technical point, reference prior theory does leave some wiggle room for models with several parameters. The user is asked to provide an ordering from least to most important. Such an ordering is not always easy to determine, especially when the ultimate goal of the study is predictive rather than inferential. For the model at hand, a partial, but not total, order on the parameters could be derived.

This situation led to the development of a new technique designed to bypass this requirement and resulting in the “Gibbs reference posterior distribution” as a way to quantify parameter uncertainty after having observed the available data in the absence of parameter ordering.

This technique relies on the novel notion of optimal compromise between potentially incompatible conditional distributions. Thus reference priors for every single parameter conditional to all others can be computed. The “Gibbs reference posterior distribution” distribution is defined as the optimal compromise between the resulting conditional reference posterior distributions. We prove that it is proper and therefore usable for Bayesian inference and prediction.

Combining the data transformation and the Gibbs reference posterior distribution produced a full solution to the surrogate model problem by naturally including uncertainty about the parameters of the model in the computation of Probability Of defect Detection curves.

Organization of the text

The dissertation is organized in three parts.

The first part, “Tools”, is a description of mathematical objects that were needed to solve the problem. It contains the “state-of-the-art” chapters.

Chapter 1 presents Kriging. Gaussian processes are defined and characterized in terms of mean function and correlation kernel. The crucial notion of stationarity is presented, and the link between stationarity of Gaussian process and stationarity of its covariance kernel is explained. The definition of mean square continuity and differentiability of Gaussian processes is then provided. Special attention is devoted to the way stationarity grounds this theory in standard differential calculus. The chapter ends with a presentation of standard covariance kernels and a short discussion of their properties.

Chapter 2 is independent from Chapter 1 and can be read first if desired. It presents the parts of reference analysis that are relevant to this thesis. It endeavors to make the formally difficult definition of a reference prior as intuitive as possible. The chapter is mainly devoted to the

one-parameter setting, because this is the setting where the theory is the most developed and yields the most interesting results. Although mainly a review of previous results, a small original contribution is added in the form of a nice uniqueness result. The final part of the chapter deals with the murkier multi-parameter setting. The approach presented is a departure from the theory established in Berger and Bernardo [1992], but it is nothing more than a formalization of the concept of “exact marginalization” prescribed by Berger et al. [2001] for Kriging models.

Chapter 3 combines notions from Chapters 1 and 2 insofar as it concerns reference analysis of Kriging models with isotropic covariance kernels. With regard to its result, it is a state-of-the-art chapter: it relates results from Berger et al. [2001]. However, it shows that Berger et al. [2001]’s proof is adequate for rough Gaussian processes only and gives an original proof for smoother processes.

The second part, “Compromise” is the shortest part – it contains a single chapter – but it plays a crucial role in this thesis.

Chapter 4 proposes a theory of compromise between possibly incompatible conditional distributions. It could be read independently from the rest of the dissertation because it does not specifically concern spatial or even Bayesian statistics. It tackles a problem that arises from the popularity of Gibbs samplers: sometimes, there is no guarantee that the conditional distributions truly define a joint distribution. However, the chapter shows that if it exists and is unique, the stationary distribution of the Markov chain associated with the Gibbs sampler can be viewed as the optimal compromise between the input conditional distributions.

The third part, “Application” details the consequences of the aforementioned theory of compromise for Objective Bayesian analysis of Kriging models with anisotropic covariance kernels.

Chapter 5 details how the theory of compromise makes it possible, in a Simple Kriging setting, not to order correlation parameters and instead use the optimal compromise between the posterior distributions derived from conditional reference priors. Numerical simulations show that this optimal compromise, the Gibbs reference posterior distribution, has good frequentist coverage properties.

Chapter 6 takes the results from Chapter 5 and extends them to the Universal Kriging setting, where the mean function of the Gaussian process is unknown. Comparable simulations are run in this setting, which is much wider than the Simple Kriging one.

Chapter 7 uses the Gibbs reference posterior distribution to tackle the industrial problem of EDF. It provides a heuristic for integrating covariance parameters out of the model thanks to the Gibbs reference posterior distribution. After a transformation family for the observation data has been defined, it makes estimating the transformation parameter easier, since all others have been integrated out. Finally, it shows how covariance parameter uncertainty can be embedded in POD curves.

Part I

Tools

Chapter 1

Kriging Overview

Abstract

The Kriging or Gaussian process regression model is introduced and its basic properties highlighted. The chapter covers Gaussian random processes, stationarity, mean square continuity and differentiability. Properties of Gaussian processes are linked to properties of their covariance kernels. Spectral analysis of Gaussian processes, beginning with Bochner's theorem, is introduced. Important examples of covariance kernels are provided, along with their spectral representations.

Résumé

Le modèle de krigeage ou régression par processus gaussiens est introduit et ses propriétés fondamentales mises en évidence. Le chapitre couvre les processus aléatoires gaussiens et les notions de stationarité, continuité en moyenne quadratique et différentiabilité en moyenne quadratique. Les propriétés des processus gaussiens sont reliées à celles de leurs noyaux de covariance. Le chapitre introduit aussi l'analyse spectrale des processus gaussiens, à commencer par le théorème de Bochner. Des exemples importants de noyaux de covariance ainsi que leur représentation spectrale sont fournis.

1.1 Introduction

This chapter provides a short introduction to Kriging from a theoretical rather than practical perspective. Its goal is to introduce tools necessary to understand the notions used in this thesis. It draws heavily on Stein [1999], Rasmussen and Williams [2006] and Bachoc [2013a]. The only prerequisites are standard probability and statistical theory.

Gaussian Stochastic Processes offer a convenient way of expressing the uncertainty about the value of some real-valued quantity on a given spatial domain \mathcal{D} [Stein, 1999] when said quantity is only observed on a finite set of points in \mathcal{D} . This is why Gaussian Process Regression is used as a supervised learning method [Rasmussen and Williams, 2006, chapter 2], although it originally appeared in the geostatistical literature [Matheron, 1960]. Throughout the dissertation, we follow the geostatistical naming convention for this model: Kriging.

1.2 Gaussian random processes

General random processes

Most random processes considered in this thesis are Gaussian. It is however useful to keep in mind the more general framework for random processes. Detailed information about this topic can be found in Billingsley [1995], Dudley [1989] or Bass [1995].

Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space. Let \mathcal{D} be a spatial or temporal domain – typically $\mathcal{D} \subset \mathbb{R}^r$ for some positive integer r . Let \mathcal{S} be a metric space – typically $\mathcal{S} = \mathbb{R}$ – and let $\mathcal{B}(\mathcal{S})$ be its Borel σ -algebra.

Definition 1.1. *A random process or random field is a mapping from the domain \mathcal{D} to the set of all random variables $(\Omega, \mathcal{F}) \rightarrow (\mathcal{S}, \mathcal{B}(\mathcal{S}))$.*

No assumption is made as to whether the domain \mathcal{D} or the metric space \mathcal{S} is discrete or continuous. Markov chains are random processes with discrete domain \mathcal{D} and discrete or continuous metric space \mathcal{S} . Poisson processes have continuous domain \mathcal{D} and discrete metric space \mathcal{S} . Wiener processes (Brownian motion) have continuous domain \mathcal{D} and continuous metric space \mathcal{S} .

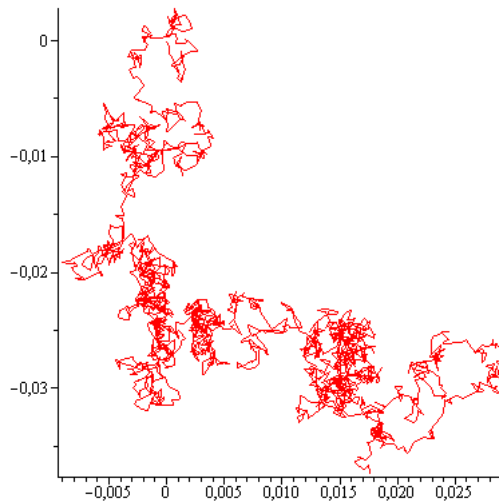


Figure 1.1 – *A trajectory of a 2-dimensional Brownian motion. Approximation through a 3000-step Markov chain where each step follows a Gaussian bi-variate distribution. Credit goes to Wikipedia contribution “Ipipipourax” who put this work in the public domain.*

An important property of Gaussian processes is that they can themselves be viewed as random variables. To see this, a few additional notions need to be defined.

Definition 1.2. *For all $\mathbf{x} \in \mathcal{D}$, define the mapping $\pi_{\mathbf{x}} : \mathbb{R}^{\mathcal{D}} \rightarrow \mathbb{R}; f \mapsto f(\mathbf{x})$. Denote by $\mathcal{B}(\mathbb{R})^{\otimes \mathcal{D}}$ the σ -algebra of $\mathbb{R}^{\mathcal{D}}$ spanned by $\{\pi_{\mathbf{x}}^{-1}(B) : \mathbf{x} \in \mathcal{D}, B \in \mathcal{B}(\mathbb{R})\}$.*

Although a random process is defined as a collection of random variables, the following proposition shows that it can also be represented as a random variable with functional value [Billingsley, 1995, chapter 7].

Proposition 1.3. *Let Y be a random process. The mapping $\tilde{Y} : \Omega \rightarrow \mathbb{R}^{\mathcal{D}}; \omega \mapsto [\mathbf{x} \mapsto Y(\mathbf{x})(\omega)]$ is $(\mathcal{F}, \mathcal{B}(\mathbb{R})^{\otimes \mathcal{D}})$ -measurable.*

Proof. Let $B \in \mathcal{B}(\mathbb{R})$. Then for all $\mathbf{x} \in \mathcal{D}$, $\tilde{Y}^{-1}(\pi_{\mathbf{x}}^{-1}(B)) = (\pi_{\mathbf{x}} \circ \tilde{Y})^{-1}(B) = Y(\mathbf{x})^{-1}(B) \in \mathcal{F}$. \square

The probability distribution of a random process is characterized by its finite-dimensional distributions [Khoshnevisan, 2002, chapter 3]:

Proposition 1.4. *Let Y_1 and Y_2 be two random processes and let \tilde{Y}_1 and \tilde{Y}_2 be the mappings $\Omega \rightarrow \mathbb{R}^{\mathcal{D}}$ defined by $\tilde{Y}_1(\omega)(\mathbf{x}) = Y_1(\mathbf{x})(\omega)$ and $\tilde{Y}_2(\omega)(\mathbf{x}) = Y_2(\mathbf{x})(\omega)$. If, for every positive integer n and every family $(\mathbf{x}^{(i)})_{i \in \llbracket 1, n \rrbracket}$ of points in \mathcal{D} , the Gaussian vectors $(Y_1(\mathbf{x}^{(1)}), \dots, Y_1(\mathbf{x}^{(n)}))^\top$ and $(Y_2(\mathbf{x}^{(1)}), \dots, Y_2(\mathbf{x}^{(n)}))^\top$ have the same probability distribution, then \tilde{Y}_1 and \tilde{Y}_2 have the same probability distribution.*

Proof. For every positive integer n , for every family $(\mathbf{x}^{(i)})_{i \in \llbracket 1, n \rrbracket}$ of elements of \mathcal{D} and every family $(B_i)_{i \in \llbracket 1, n \rrbracket}$ of elements of $\mathcal{B}(\mathbb{R})$,

$$\begin{aligned} \tilde{Y}_1^{-1}(\pi_{\mathbf{x}^{(1)}}^{-1}(B_1) \cap \dots \cap \pi_{\mathbf{x}^{(n)}}^{-1}(B_n)) &= \tilde{Y}_1^{-1}(\pi_{\mathbf{x}^{(1)}}^{-1}(B_1)) \cap \dots \cap \tilde{Y}_1^{-1}(\pi_{\mathbf{x}^{(n)}}^{-1}(B_n)) \\ &= (\pi_{\mathbf{x}^{(1)}} \circ \tilde{Y}_1)^{-1}(B_1) \cap \dots \cap (\pi_{\mathbf{x}^{(n)}} \circ \tilde{Y}_1)^{-1}(B_n) \\ &= Y_1(\mathbf{x}^{(1)})^{-1}(B_1) \cap \dots \cap Y_1(\mathbf{x}^{(n)})^{-1}(B_n). \end{aligned} \quad (1.1)$$

Similarly,

$$\tilde{Y}_2^{-1}(\pi_{\mathbf{x}^{(1)}}^{-1}(B_1) \cap \dots \cap \pi_{\mathbf{x}^{(n)}}^{-1}(B_n)) = Y_2(\mathbf{x}^{(1)})^{-1}(B_1) \cap \dots \cap Y_2(\mathbf{x}^{(n)})^{-1}(B_n). \quad (1.2)$$

Using the assumption,

$$\mathbb{P}\left(Y_1(\mathbf{x}^{(1)})^{-1}(B_1) \cap \dots \cap Y_1(\mathbf{x}^{(n)})^{-1}(B_n)\right) = \mathbb{P}\left(Y_2(\mathbf{x}^{(1)})^{-1}(B_1) \cap \dots \cap Y_2(\mathbf{x}^{(n)})^{-1}(B_n)\right). \quad (1.3)$$

Combining these equations, we obtain

$$\mathbb{P}\left(\tilde{Y}_1^{-1}(\pi_{\mathbf{x}^{(1)}}^{-1}(B_1) \cap \dots \cap \pi_{\mathbf{x}^{(n)}}^{-1}(B_n))\right) = \mathbb{P}\left(\tilde{Y}_2^{-1}(\pi_{\mathbf{x}^{(1)}}^{-1}(B_1) \cap \dots \cap \pi_{\mathbf{x}^{(n)}}^{-1}(B_n))\right). \quad (1.4)$$

Consider the following set:

$$\left\{ \pi_{\mathbf{x}^{(1)}}^{-1}(B_1) \cap \dots \cap \pi_{\mathbf{x}^{(n)}}^{-1}(B_n) : n \in \mathbb{Z}_+, (\mathbf{x}^{(i)})_{i \in \llbracket 1, n \rrbracket} \in \mathcal{D}^n, (B_i)_{i \in \llbracket 1, n \rrbracket} \in \mathcal{B}(\mathbb{R})^n \right\}.$$

It is a π -system which includes $\{\pi_{\mathbf{x}}^{-1}(B) : \mathbf{x} \in \mathcal{D}, B \in \mathcal{B}(\mathbb{R})\}$, the set spanning $\mathcal{B}(\mathbb{R})^{\otimes \mathcal{D}}$. The distributions of \tilde{Y}_1 and \tilde{Y}_2 agree on this π -system and can therefore only be equal. \square

From this point onwards, for any random process Y , we will abuse notations and also denote by Y the mapping \tilde{Y} defined in Proposition 1.3.

Gaussian random processes

The theory of Gaussian vectors, i.e. of the multivariate Normal distribution, is a prerequisite of the following. Please refer to Appendix B.1 of Santner et al. [2003] or to Appendix A.2 of Rasmussen and Williams [2006].

Gaussian random processes or Gaussian processes for short are a particular type of random processes. The most prominent members of this class are Wiener processes. For details, see for example Khoshnevisan [2002, chapter 5].

Definition 1.5. *A Gaussian process is a random process Y with metric space $\mathcal{S} = \mathbb{R}$ such that for every positive integer n and every family $(\mathbf{x}^{(i)})_{i \in [1, n]}$, $(Y(\mathbf{x}^{(1)}), \dots, Y(\mathbf{x}^{(n)}))^\top$ is a Gaussian vector.*

The class of Gaussian processes is, among all classes of random processes, the one that has been most studied. Many nice properties are available, and we only use some of them. For a larger view, please refer to Adler [1990] or Stein [1999].

Conditional Gaussian processes

Recall the following result from Gaussian vector theory:

Theorem 1.6. *Let $(\mathbf{V}_1^\top, \mathbf{V}_2^\top)^\top$ be a Gaussian vector. Assume that \mathbf{V}_1 is nondegenerate. Then its covariance matrix $\text{Var}[\mathbf{V}_1]$ is nonsingular and the distribution of \mathbf{V}_2 conditionally to \mathbf{V}_1 is Gaussian with mean vector*

$$\mathbb{E}[\mathbf{V}_2 | \mathbf{V}_1] = \mathbb{E}[\mathbf{V}_2] - \text{Cov}[\mathbf{V}_2, \mathbf{V}_1] \text{Var}[\mathbf{V}_1]^{-1} (\mathbf{V}_1 - \mathbb{E}[\mathbf{V}_1]) \quad (1.5)$$

and covariance matrix

$$\text{Var}[\mathbf{V}_2 | \mathbf{V}_1] = \text{Cov}[\mathbf{V}_2, \mathbf{V}_1] \text{Var}[\mathbf{V}_1]^{-1} \text{Cov}[\mathbf{V}_1, \mathbf{V}_2]. \quad (1.6)$$

This theorem is of great practical importance in the study of Gaussian processes. It can be generalized as follows:

Corollary 1.7. *Let $(\mathbf{V}_1^\top, \mathbf{V}_2^\top)^\top$ be a Gaussian vector. If n is the size of \mathbf{V} and w is the rank of $\text{Var}[\mathbf{V}_1]$, let \mathbf{W} be an $n \times w$ matrix representing an isometry from the subspace of \mathbb{R}^n spanned by $\text{Var}[\mathbf{V}_1]$ to \mathbb{R}^w . Then the distribution of \mathbf{V}_2 conditionally to \mathbf{V}_1 is Gaussian with mean vector*

$$\mathbb{E}[\mathbf{V}_2 | \mathbf{V}_1] = \mathbb{E}[\mathbf{V}_2] - \text{Cov}[\mathbf{V}_2, \mathbf{V}_1] \mathbf{W} \left(\mathbf{W}^\top \text{Var}[\mathbf{V}_1] \mathbf{W} \right)^{-1} \mathbf{W}^\top (\mathbf{V}_1 - \mathbb{E}[\mathbf{V}_1]) \quad (1.7)$$

and covariance matrix

$$\text{Var}[\mathbf{V}_2 | \mathbf{V}_1] = \text{Cov}[\mathbf{V}_2, \mathbf{V}_1] \mathbf{W} \left(\mathbf{W}^\top \text{Var}[\mathbf{V}_1] \mathbf{W} \right)^{-1} \mathbf{W}^\top \text{Cov}[\mathbf{V}_1, \mathbf{V}_2]. \quad (1.8)$$

Proof. The $n \times n$ matrix $\mathbf{W}\mathbf{W}^\top$ represents an orthogonal projection on the subspace of \mathbb{R}^n spanned by $\text{Var}[\mathbf{V}_1]$. So $\mathbf{V}_1 = \mathbf{W}\mathbf{W}^\top \mathbf{V}_1$, which implies that the distribution of \mathbf{V}_2 conditionally to \mathbf{V}_1 is the distribution of \mathbf{V}_2 conditionally to $\mathbf{W}^\top \mathbf{V}_1$.

Since $(\mathbf{V}_1^\top, \mathbf{V}_2^\top)^\top$ is a Gaussian vector, any linear transformation remains a Gaussian vector. In particular, $((\mathbf{W}\mathbf{V}_1)^\top, \mathbf{V}_2^\top)^\top$ is a Gaussian vector. Moreover, $\mathbf{W}^\top \mathbf{V}_1$ is nondegenerate so Theorem 1.6 is applicable.

The distribution of \mathbf{V}_2 conditionnally to \mathbf{V}_1 (or equivalently of $\mathbf{W}^\top \mathbf{V}_1$) is Gaussian with mean vector

$$\begin{aligned} \mathbb{E}[\mathbf{V}_2|\mathbf{V}_1] &= \mathbb{E}[\mathbf{V}_2] - \text{Cov}[\mathbf{V}_2, \mathbf{W}\mathbf{V}_1] \text{Var}[\mathbf{W}^\top \mathbf{V}_1 \mathbf{W}]^{-1} \text{Cov}[\mathbf{W}\mathbf{V}_1, \mathbf{V}_2] \\ &= \mathbb{E}[\mathbf{V}_2] - \text{Cov}[\mathbf{V}_2, \mathbf{V}_1] \mathbf{W} \left(\mathbf{W}^\top \text{Var}[\mathbf{V}_1] \mathbf{W} \right)^{-1} \mathbf{W}^\top (\mathbf{V}_1 - \mathbb{E}[\mathbf{V}_1]) \end{aligned} \quad (1.9)$$

and with covariance matrix

$$\begin{aligned} \text{Var}[\mathbf{V}_2|\mathbf{V}_1] &= \text{Cov}[\mathbf{V}_2, \mathbf{W}\mathbf{V}_1] \text{Var}[\mathbf{W}\mathbf{V}_1]^{-1} \mathbf{W}^\top \text{Cov}[\mathbf{W}\mathbf{V}_1, \mathbf{V}_2] \\ &= \text{Cov}[\mathbf{V}_2, \mathbf{V}_1] \mathbf{W} \left(\mathbf{W}^\top \text{Var}[\mathbf{V}_1] \mathbf{W} \right)^{-1} \mathbf{W}^\top \text{Cov}[\mathbf{V}_1, \mathbf{V}_2]. \end{aligned} \quad (1.10)$$

□

With this result, we can consider conditional Gaussian processes [Stein, 1999, chapter 1].

Definition 1.8. *Let Y be a Gaussian process. Let n be a positive integer and $(\mathbf{x}^{(i)})_{i \in [1, n]}$ be a family of n points in \mathcal{D} . A Gaussian process Y_{cond} such that for every positive integer n' and every family $(\mathbf{x}^{(i')})_{i' \in [1, n']}$ of points in \mathcal{D} , the Gaussian vector $(Y_{\text{cond}}(\mathbf{x}^{(1')}), \dots, Y_{\text{cond}}(\mathbf{x}^{(n')}))^\top$ has the same distribution as the Gaussian vector $(Y(\mathbf{x}^{(1')}), \dots, Y(\mathbf{x}^{(n')}))^\top$ conditionally to $Y(\mathbf{x}^{(1)}), \dots, Y(\mathbf{x}^{(n)})$ is called a version of the conditional process Y knowing $Y(\mathbf{x}^{(1)}), \dots, Y(\mathbf{x}^{(n)})$.*

It is easy to see that all versions of the conditional process have the same distribution.

This notion of conditional Gaussian process is at the heart of the Kriging procedure.

Definition 1.9. *Let f be an unknown mapping $\mathcal{D} \rightarrow \mathbb{R}$. Assume that for some positive integer n , there exists a family of n points in \mathcal{D} $(\mathbf{x}^{(i)})_{i \in [1, n]}$ on which f was observed, so $f(\mathbf{x}^{(1)}), \dots, f(\mathbf{x}^{(n)})$ are known. A Kriging model is a version of the conditional Gaussian process Y knowing $Y(\mathbf{x}^{(1)}) = f(\mathbf{x}^{(1)}), \dots, Y(\mathbf{x}^{(n)}) = f(\mathbf{x}^{(n)})$.*

Since all versions of the conditional Gaussian process have the same distribution, distinguishing them makes no sense from a statistical perspective. In the following, we abuse denonimiations by referring to any version of the conditional process as *the* conditional process.

In this construction, the distribution of the conditional Gaussian process Y is supposed to represent the uncertainty on f . It does *not* mean that f is assumed to be a realization of a Gaussian process, which would be a meaningless statement.

Characterizing a Gaussian process

Just like mean vector and covariance matrix characterize the distribution of a Gaussian vector, mean function and covariance function (or kernel) characterize a Gaussian process [Bachoc, 2013a, chapter 2].

Definition 1.10. Let Y be a Gaussian process. The mapping $m : \mathcal{D} \rightarrow \mathbb{R}$ defined by $\mathbf{x} \mapsto \mathbb{E}[Y(\mathbf{x})]$ is called the mean function of Y .

Definition 1.11. Let Y be a Gaussian process. The mapping $K : \mathcal{D} \times \mathcal{D} \rightarrow \mathbb{R}$ defined by $(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \mapsto \text{Cov}[Y(\mathbf{x}^{(1)}), Y(\mathbf{x}^{(2)})]$ is called the covariance function or covariance kernel of Y .

Proposition 1.12. Let Y_1 and Y_2 be Gaussian processes with the same mean function and covariance kernel. Then they have the same distribution.

Proof. For any positive integer n and any family $(\mathbf{x}^{(i)})_{i \in [1, n]}$ of n points in \mathcal{D} , the random vectors $(Y_1(\mathbf{x}^{(1)}), \dots, Y_1(\mathbf{x}^{(n)}))^\top$ and $(Y_2(\mathbf{x}^{(1)}), \dots, Y_2(\mathbf{x}^{(n)}))^\top$ are Gaussian vectors with the same mean vector and covariance matrix, so their distributions are equal. The result then follows from Proposition 1.4. \square

If mean function and covariance kernel are enough to characterize the distribution of a Gaussian process, it raises the question of which functions are admissible as mean function and covariance kernel. That is, for which mean function and covariance kernel candidates is it possible to construct a Gaussian process that admits them as mean function and covariance kernel respectively?

Provided a Gaussian process with null mean function can be constructed, a Gaussian process with any mean function can be obtained simply by adding said function. So there is no condition on mean functions.

The picture is different for covariance kernels.

Definition 1.13. A mapping $K : \mathcal{D} \times \mathcal{D} \rightarrow \mathbb{R}$ is called positive definite if for every positive integer n and any family $(\mathbf{x}^{(i)})_{i \in [1, n]}$ of n points in \mathcal{D} , the $n \times n$ matrix with (i, j) -th element $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ is positive semi-definite.

A positive definite mapping is sometimes called a kernel, hence the name ‘‘covariance kernel’’ for ‘‘covariance function’’. See the study of the positive-definiteness of bivariate mappings conducted in Schölkopf and Smola [2012][chapter 13].

Proposition 1.14. For any mapping $m : \mathcal{D} \rightarrow \mathbb{R}$ and any mapping $K : \mathcal{D} \times \mathcal{D} \rightarrow \mathbb{R}$, there exists a Gaussian process Y with mean function m and covariance kernel K if and only if K is positive definite.

In the context of the proposition above, the phrase ‘‘there exists a Gaussian process’’ means that there exists a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ and a corresponding random process $Y : \mathcal{D} \mapsto \mathbb{R}$ which is a Gaussian process. In all preceding and following results, the existence of $(\Omega, \mathcal{F}, \mathbb{P})$ is taken for granted. This is the only one in which this is not the case. Obviously, for K to be the covariance kernel of some Gaussian process, it needs to be positive definite, so the condition given in the proposition is necessary. It is its sufficiency – i.e. the existence of $(\Omega, \mathcal{F}, \mathbb{P})$ and Y – that is difficult to show. It requires Kolmogorov’s existence theorem, of which Billingsley [1995, chapter 7] provides two different proofs.

Stationarity

General Gaussian processes can be hard to handle. Practically speaking, most Kriging models used in the literature consider a particular class of Gaussian processes: stationary Gaussian processes [Rasmussen and Williams, 2006, chapter 4]. To properly define them, we need to assume that \mathcal{D} is a subset of a real affine space. Let \mathcal{V} be the corresponding real vector space.

Definition 1.15. *A Gaussian process Y is said to be stationary if for any positive integer n , any family $(\mathbf{x}^{(i)})_{i \in [1, n]}$ of n points in \mathcal{D} and any vector $\mathbf{h} \in \mathcal{V}$, the following holds: if for every integer $i \in [1, n]$, $\mathbf{x}^{(i)} + \mathbf{h} \in \mathcal{D}$, then the Gaussian vectors $(Y(\mathbf{x}^{(1)} + \mathbf{h}), \dots, Y(\mathbf{x}^{(n)} + \mathbf{h}))^\top$ and $(Y(\mathbf{x}^{(1)}), \dots, Y(\mathbf{x}^{(n)}))^\top$ share the same distribution.*

Stationarity can be characterized in terms of mean function and covariance kernel [Rasmussen and Williams, 2006, chapter 4].

Definition 1.16. *A positive definite mapping $K : \mathcal{D} \times \mathcal{D} \rightarrow \mathbb{R}$ is said to be stationary if there exists a mapping $\tilde{K} : \mathcal{V} \rightarrow \mathbb{R}$ such that for every $(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \in \mathcal{D}^2$, $K(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) = \tilde{K}(\mathbf{x}^{(1)} - \mathbf{x}^{(2)})$.*

Proposition 1.17. *A Gaussian process is stationary if and only if its mean function is constant and its covariance kernel stationary.*

Proof. Let Y be a Gaussian process with mean function m and covariance kernel K .

Assume that m is constant and that K is stationary. Then there exists a mapping $\tilde{K} : \mathcal{V} \rightarrow \mathbb{R}$ such that for every $(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \in \mathcal{D}^2$, $K(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) = \tilde{K}(\mathbf{x}^{(1)} - \mathbf{x}^{(2)})$.

Let n be a positive integer and let $(\mathbf{x}^{(i)})_{i \in [1, n]}$ be a family of n points in \mathcal{D} . Let \mathbf{h} be a vector of \mathcal{V} such that for every integer $i \in [1, n]$, $\mathbf{x}^{(i)} + \mathbf{h} \in \mathcal{D}$. Then the Gaussian vector $(Y(\mathbf{x}^{(1)} + \mathbf{h}), \dots, Y(\mathbf{x}^{(n)} + \mathbf{h}))^\top$ has constant mean vector of value $m(\mathbf{x}^{(1)})$ and its covariance matrix has (i, j) -th element $K(\mathbf{x}^{(i)} + \mathbf{h}, \mathbf{x}^{(j)} + \mathbf{h}) = \tilde{K}(\mathbf{x}^{(i)} + \mathbf{h} - (\mathbf{x}^{(j)} + \mathbf{h})) = \tilde{K}(\mathbf{x}^{(i)} - \mathbf{x}^{(j)}) = K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$. So its mean vector and covariance matrix are respectively equal to the mean vector and covariance matrix of the Gaussian vector $(Y(\mathbf{x}^{(1)}), \dots, Y(\mathbf{x}^{(n)}))^\top$. Therefore Y is stationary.

Now assume that Y is stationary. Let $\mathbf{x}^{(1)} \in \mathcal{D}$. Then, for all $\mathbf{x}^{(2)} \in \mathcal{D}$, $Y(\mathbf{x}^{(2)}) = Y(\mathbf{x}^{(1)} + (\mathbf{x}^{(2)} - \mathbf{x}^{(1)}))$ has the same distribution as $Y(\mathbf{x}^{(1)})$. In particular, it has the same mean. Therefore the mean function m is constant.

Let $(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \in \mathcal{D}^2$ and $(\mathbf{x}^{(1)'}, \mathbf{x}^{(2)'}) \in \mathcal{D}^2$ such that $\mathbf{x}^{(1)} - \mathbf{x}^{(2)} = \mathbf{x}^{(1)'} - \mathbf{x}^{(2)'}$. Then the Gaussian vector $(Y(\mathbf{x}^{(1)'}, Y(\mathbf{x}^{(2)'}))^\top = (Y(\mathbf{x}^{(1)} + (\mathbf{x}^{(1)'} - \mathbf{x}^{(1)})), Y(\mathbf{x}^{(2)} + (\mathbf{x}^{(2)'} - \mathbf{x}^{(2)})))^\top = (Y(\mathbf{x}^{(1)} + (\mathbf{x}^{(1)'} - \mathbf{x}^{(1)})), Y(\mathbf{x}^{(2)} + (\mathbf{x}^{(1)'} - \mathbf{x}^{(1)})))^\top$ has the same distribution as the Gaussian vector $(Y(\mathbf{x}^{(1)}), Y(\mathbf{x}^{(2)}))^\top$. In particular, the covariance $K(\mathbf{x}^{(1)'}, \mathbf{x}^{(2)'})$ between $Y(\mathbf{x}^{(1)'})$ and $Y(\mathbf{x}^{(2)'})$ is equal to the covariance $K(\mathbf{x}^{(1)}, \mathbf{x}^{(2)})$ between $Y(\mathbf{x}^{(1)})$ and $Y(\mathbf{x}^{(2)})$. So the mapping \tilde{K} can be defined as follows: for every $v \in \mathcal{V}$, if there exists a pair $(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \in \mathcal{D}^2$ such that $\mathbf{x}^{(1)} - \mathbf{x}^{(2)} = v$, then $\tilde{K}(v) := K(\mathbf{x}^{(1)}, \mathbf{x}^{(2)})$. Else $\tilde{K}(v) := 0$.

□

From this point onwards, we abuse denominations and call $\tilde{K} : \mathcal{V} \rightarrow \mathbb{R}$ the covariance kernel instead of $K : \mathcal{D} \times \mathcal{D} \rightarrow \mathbb{R}$. This is a twofold abuse because \tilde{K} is not necessarily unique, as the proof above shows. It is of no matter, because \tilde{K} is uniquely defined at all points that can be written $\mathbf{x}^{(1)} - \mathbf{x}^{(2)}$ with $(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \in \mathcal{D}^2$ and thus K can be computed from \tilde{K} .

Since positive definite mappings are covariance kernels, this abuse leads to an extension of the definition of positive definiteness to mappings $\mathbb{R}^r \rightarrow \mathbb{R}$.

Definition 1.18. *A mapping $f : \mathbb{R}^r \rightarrow \mathbb{R}$ is called positive definite if there exists a positive definite mapping $K : \mathbb{R}^r \times \mathbb{R}^r \rightarrow \mathbb{R}$ such that for all $(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \in \mathbb{R}^r \times \mathbb{R}^r$, $f(\mathbf{x}^{(1)} - \mathbf{x}^{(2)}) = K(\mathbf{x}^{(1)}, \mathbf{x}^{(2)})$.*

1.3 Mean square continuity and differentiability of Gaussian processes

The L^2 norm allows a natural definition of continuity and differentiability for Gaussian processes that is linked to the continuity and stationarity of its covariance kernel.

As before, we assume that \mathcal{D} is a subset of a real affine space with corresponding vector space \mathcal{V} . Let us further assume that $\mathcal{V} = \mathbb{R}^r$ for a given $r \in \mathbb{N}$. We endow the affine space with the metric induced by the Euclidean norm on \mathbb{R}^r and assume that \mathcal{D} is an open subset.

Mean square continuity and differentiability

Let us consider mean square properties of Gaussian processes [Stein, 1999, chapter 2].

Definition 1.19. *A Gaussian process Y is mean square continuous on \mathcal{D} if for any $\mathbf{x}^{(0)} \in \mathcal{D}$ and for all $\epsilon > 0$, there exists $\delta > 0$ such that for every $\mathbf{x} \in \mathcal{D}$ such that $\|\mathbf{x} - \mathbf{x}^{(0)}\| \leq \delta$,*

$$\mathbb{E} \left[\left(Y(\mathbf{x}) - Y(\mathbf{x}^{(0)}) \right)^2 \right] \leq \epsilon^2 \quad (1.11)$$

Remark. For a stationary process, being mean square continuous is equivalent to being mean square uniformly continuous. This means that in the above definition, δ does not depend on $\mathbf{x}^{(0)}$.

Definition 1.20. *A Gaussian process Y is mean square continuously differentiable on \mathcal{D} if there exist r mean square continuous Gaussian processes $\partial_{\mathbf{e}_1} Y, \dots, \partial_{\mathbf{e}_r} Y$ such that, defining $(\mathbf{e}_k)_{k \in \llbracket 1, r \rrbracket}$ as the canonical basis of \mathbb{R}^r , for any integer $k \in \llbracket 1, r \rrbracket$, the following property holds. For any $\mathbf{x}^{(0)} \in \mathcal{D}$ and for all $\epsilon > 0$, there exists $\delta > 0$ such that for all $0 < h \leq \delta$,*

$$\mathbb{E} \left[\left(\frac{Y(\mathbf{x}^{(0)} + h\mathbf{e}_k) - Y(\mathbf{x}^{(0)})}{h} - \partial_{\mathbf{e}_k} Y(\mathbf{x}^{(0)}) \right)^2 \right] \leq \epsilon^2 \quad (1.12)$$

In the definition above, for every $k \in \llbracket 1, r \rrbracket$, the ‘‘partial derivative’’ $\partial_{\mathbf{e}_k} Y$ is unique in the sense that any other Gaussian process fitting the requirement would necessarily be a modification of $\partial_{\mathbf{e}_k} Y$.

Remark. If Y is a stationary Gaussian process that is continuously differentiable, then all $\partial_{\mathbf{e}_k} Y$ are necessarily stationary Gaussian processes.

The following result provides an alternate view of continuous mean square differentiability [Bachoc, 2013a, chapter 2].

Proposition 1.21. *A stationary Gaussian process Y is mean square continuously differentiable on \mathcal{D} if, and only if, $Y : \mathcal{D} \rightarrow L^2(\Omega)$ is continuously Fréchet differentiable on \mathcal{D} .*

Proof. If a stationary Gaussian process Y is Fréchet continuously differentiable on \mathcal{D} , then for all $\mathbf{x}^{(0)} \in \mathcal{D}$ and every integer $k \in \llbracket 1, r \rrbracket$, consider the value taken by the Fréchet derivative $DY(\mathbf{x}^{(0)})$ at the vector \mathbf{e}_k .

Equation (1.12) holds with $\partial_{\mathbf{e}_k} Y(\mathbf{x}^{(0)}) := DY(\mathbf{x}^{(0)})(\mathbf{e}_k)$ by definition of $DY(\mathbf{x}^{(0)})(\mathbf{e}_k)$. Moreover, the Fréchet derivative $DY(\mathbf{x}^{(0)})(\mathbf{e}_k)$ is continuous. To prove that Y is mean square continuously differentiable, we only need to show that the mapping $\mathbf{x}^{(0)} \in \mathcal{D} \mapsto DY(\mathbf{x}^{(0)})(\mathbf{e}_k)$ is a Gaussian process. This follows from the fact that for any positive integer n and any points $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)} \in \mathcal{D}$, $(DY(\mathbf{x}^{(1)})(\mathbf{e}_k), \dots, DY(\mathbf{x}^{(n)})(\mathbf{e}_k))^\top$ is an L^2 limit of Gaussian vectors and is therefore a Gaussian vector.

Conversely, if Y is a stationary Gaussian process that is mean square differentiable, define for every point $\mathbf{x}^{(0)} \in \mathcal{D}$ the linear application $DY(\mathbf{x}^{(0)}) : \mathbb{R}^r \rightarrow L^2(\Omega)$ such that for every $k \in \llbracket 1, r \rrbracket$, $DY(\mathbf{x}^{(0)})(\mathbf{e}_k) = \partial_{\mathbf{e}_k} Y(\mathbf{x}^{(0)})$.

For every $\mathbf{x}^{(0)} \in \mathcal{D}$ and every $\mathbf{u} \in \mathbb{R}^r$ and $\mathbf{v} \in \mathbb{R}^r$ such that $\mathbf{x}^{(0)} + \mathbf{u} \in \mathcal{D}$ and $\mathbf{x}^{(0)} + \mathbf{u} + \mathbf{v} \in \mathcal{D}$, because $DY(\mathbf{x}^{(0)})$ is a linear mapping, we have

$$Y(\mathbf{x}^{(0)} + \mathbf{u} + \mathbf{v}) - Y(\mathbf{x}^{(0)}) - DY(\mathbf{x}^{(0)})(\mathbf{u} + \mathbf{v}) \quad (1.13)$$

$$\begin{aligned} &= Y(\mathbf{x}^{(0)} + \mathbf{u} + \mathbf{v}) - Y(\mathbf{x}^{(0)} + \mathbf{u}) - DY(\mathbf{x}^{(0)} + \mathbf{u})(\mathbf{v}) \\ &\quad + Y(\mathbf{x}^{(0)} + \mathbf{u}) - Y(\mathbf{x}^{(0)}) - DY(\mathbf{x}^{(0)})(\mathbf{u}) \\ &\quad + DY(\mathbf{x}^{(0)} + \mathbf{u})(\mathbf{v}) - DY(\mathbf{x}^{(0)})(\mathbf{v}). \end{aligned} \quad (1.14)$$

Because \mathcal{D} is open, for every $\mathbf{x}^{(0)} \in \mathcal{D}$, there exists $t > 0$ such that for every $\mathbf{h} \in \mathbb{R}^r$ such that $\|\mathbf{h}\| < t$, $\mathbf{x}^{(0)} + \mathbf{h} \in \mathcal{D}$. Let $h_1 \mathbf{e}_1 + \dots + h_r \mathbf{e}_r$ ($h_1, \dots, h_r \in \mathbb{R}$) be the unique decomposition of \mathbf{h} in the canonical basis of \mathbb{R}^r . Applying the result above, we have

$$\begin{aligned} &Y(\mathbf{x}^{(0)} + \mathbf{h}) - Y(\mathbf{x}^{(0)}) - DY(\mathbf{x}^{(0)})(\mathbf{u}) \\ &= Y\left(\mathbf{x}^{(0)} + \sum_{k=1}^r h_k \mathbf{e}_k\right) - Y(\mathbf{x}^{(0)}) - DY(\mathbf{x}^{(0)})\left(\sum_{k=1}^r h_k \mathbf{e}_k\right) \\ &= \sum_{k=1}^r Y\left(\mathbf{x}^{(0)} + \sum_{l=1}^k h_l \mathbf{e}_l\right) - Y\left(\mathbf{x}^{(0)} + \sum_{l=1}^{k-1} h_l \mathbf{e}_l\right) - DY\left(\mathbf{x}^{(0)} + \sum_{l=1}^{k-1} h_l \mathbf{e}_l\right)(h_k \mathbf{e}_k) \\ &\quad + \sum_{k=1}^r DY\left(\mathbf{x}^{(0)} + \sum_{l=1}^{k-1} h_l \mathbf{e}_l\right)(h_k \mathbf{e}_k) - DY(\mathbf{x}^{(0)})(h_k \mathbf{e}_k) \\ &= \sum_{k=1}^r Y\left(\mathbf{x}^{(0)} + \sum_{l=1}^k h_l \mathbf{e}_l\right) - Y\left(\mathbf{x}^{(0)} + \sum_{l=1}^{k-1} h_l \mathbf{e}_l\right) - h_k \partial_{\mathbf{e}_k} Y\left(\mathbf{x}^{(0)} + \sum_{l=1}^{k-1} h_l \mathbf{e}_l\right) \\ &\quad + \sum_{k=1}^r h_k \partial_{\mathbf{e}_k} Y\left(\mathbf{x}^{(0)} + \sum_{l=1}^{k-1} h_l \mathbf{e}_l\right) - h_k \partial_{\mathbf{e}_k} Y(\mathbf{x}^{(0)}) \end{aligned} \quad (1.15)$$

In the following, $\|\cdot\|$ denotes the Euclidean norm if applied to a vector of \mathbb{R}^r and the L^2 norm if applied to a random variable.

For every $k \in \llbracket 1, r \rrbracket$, due to stationarity,

$$\begin{aligned} & \left\| Y \left(\mathbf{x}^{(0)} + \sum_{l=1}^k h_l \mathbf{e}_l \right) - Y \left(\mathbf{x}^{(0)} + \sum_{l=1}^{k-1} h_l \mathbf{e}_l \right) - h_k \partial_{\mathbf{e}_k} Y \left(\mathbf{x}^{(0)} + \sum_{l=1}^{k-1} h_l \mathbf{e}_l \right) \right\| \\ &= \left\| Y \left(\mathbf{x}^{(0)} + h_k \mathbf{e}_k \right) - Y \left(\mathbf{x}^{(0)} \right) - h_k \partial_{\mathbf{e}_k} Y \left(\mathbf{x}^{(0)} \right) \right\|. \end{aligned} \quad (1.16)$$

So combining both equations above yields

$$\begin{aligned} & \frac{1}{\|\mathbf{h}\|} \left\| Y(\mathbf{x}^{(0)} + \mathbf{h}) - Y(\mathbf{x}^{(0)}) - DY(\mathbf{x}^{(0)})(\mathbf{u}) \right\| \\ & \leq \sum_{k=1}^r \frac{|h_k|}{\|\mathbf{h}\|} \left\| \frac{Y(\mathbf{x}^{(0)} + h_k \mathbf{e}_k) - Y(\mathbf{x}^{(0)})}{h_k} - \partial_{\mathbf{e}_k} Y(\mathbf{x}^{(0)}) \right\| \\ & \quad + \sum_{k=1}^r \frac{|h_k|}{\|\mathbf{h}\|} \left\| \partial_{\mathbf{e}_k} Y \left(\mathbf{x}^{(0)} + \sum_{l=1}^{k-1} h_l \mathbf{e}_l \right) - \partial_{\mathbf{e}_k} Y(\mathbf{x}^{(0)}) \right\|. \end{aligned} \quad (1.17)$$

This leads to

$$\begin{aligned} & \frac{1}{\|\mathbf{h}\|} \left\| Y(\mathbf{x}^{(0)} + \mathbf{h}) - Y(\mathbf{x}^{(0)}) - DY(\mathbf{x}^{(0)})(\mathbf{u}) \right\| \\ & \leq \sum_{k=1}^r \left\| \frac{Y(\mathbf{x}^{(0)} + h_k \mathbf{e}_k) - Y(\mathbf{x}^{(0)})}{h_k} - \partial_{\mathbf{e}_k} Y(\mathbf{x}^{(0)}) \right\| \\ & \quad + \sum_{k=1}^r \left\| \partial_{\mathbf{e}_k} Y \left(\mathbf{x}^{(0)} + \sum_{l=1}^{k-1} h_l \mathbf{e}_l \right) - \partial_{\mathbf{e}_k} Y(\mathbf{x}^{(0)}) \right\|. \end{aligned} \quad (1.18)$$

Because Y is continuously mean square differentiable, when $\|\mathbf{h}\| \rightarrow 0$ (and therefore every $h_k \rightarrow 0$), the right side of the inequality has null limit. The left side has therefore the same limit and Y is Fréchet differentiable at $\mathbf{x}^{(0)}$. Continuity of the Fréchet derivative follows once again from stationarity. \square

This Proposition roots the notion of mean square continuous differentiability for stationary Gaussian processes in standard differential calculus. We may therefore write that a stationary Gaussian process is of class C^0 if it is mean square continuous and of class C^1 if it is mean square continuously differentiable. Moreover, the class C^n is well defined for every nonnegative integer n .

In particular, standard differential calculus theory yields the following characterization of the class C^n .

Proposition 1.22. *A stationary Gaussian process Y is of class C^n on \mathcal{D} (i.e. n times mean square continuously differentiable) if, and only if, all partial derivatives $\partial_{\mathbf{e}_{k_1}} \partial_{\mathbf{e}_{k_2}} \dots \partial_{\mathbf{e}_{k_n}} Y$ ($k_1, k_2, \dots, k_n \in \llbracket 1, r \rrbracket$) exist and are mean square continuous on \mathcal{D} .*

Mean square continuity of stationary Gaussian processes can be characterized in terms of their covariance kernel, as the two following results show [Bachoc, 2013a, chapter 2].

Proposition 1.23. *A stationary Gaussian process Y is mean square continuous on \mathcal{D} if, and only if, its covariance kernel $K : \mathbb{R}^r \rightarrow [0, +\infty)$ is continuous at $\mathbf{0}$.*

Proof. Without loss of generality, we may assume the mean function of Y to be null. Then

$$\begin{aligned} \mathbb{E} \left[\left(Y(\mathbf{x}) - Y(\mathbf{x}^{(0)}) \right)^2 \right] &= \mathbb{E} [Y(\mathbf{x})^2] + \mathbb{E} [Y(\mathbf{x}^{(0)})^2] - 2\mathbb{E} [Y(\mathbf{x})Y(\mathbf{x}^{(0)})] \\ &= 2K(\mathbf{0}) - 2K(\mathbf{x} - \mathbf{x}^{(0)}). \end{aligned} \quad (1.19)$$

This equality implies the result. \square

Continuity of the covariance kernel at $\mathbf{0}$ is determinant, because it actually implies continuity everywhere (that matters). To make this notion precise, define $\mathbb{R}_{\mathcal{D}}^r := \{\mathbf{t} \in \mathbb{R}^r : (\exists \mathbf{x}^{(0)} \in \mathcal{D}) \mathbf{x}^{(0)} + \mathbf{t} \in \mathcal{D}\}$. $\mathbb{R}_{\mathcal{D}}^r$ is a non-empty ($\mathbf{0} \in \mathbb{R}_{\mathcal{D}}^r$) open set because \mathcal{D} is open.

Proposition 1.24. *A stationary Gaussian process Y is mean square continuous on \mathcal{D} if, and only if, its covariance kernel $K : \mathbb{R}^r \rightarrow [0, +\infty)$ is continuous on $\mathbb{R}_{\mathcal{D}}^r$.*

Proof. Because of Proposition 1.23, we only need to prove that if a stationary Gaussian process Y is mean square continuous, then its covariance kernel is continuous on $\mathbb{R}_{\mathcal{D}}^r$.

Let $\mathbf{t} \in \mathbb{R}_{\mathcal{D}}^r$. Then, because \mathcal{D} is open, for all $\mathbf{h} \in \mathbb{R}^r$ such that $\|\mathbf{h}\|$ is small enough, there exists $\mathbf{x}^{(0)} \in \mathcal{D}$ such that $\mathbf{x}^{(0)} + \mathbf{t} + \mathbf{h} \in \mathcal{D}$ and $\mathbf{x}^{(0)} + \mathbf{h} \in \mathcal{D}$.

$$\begin{aligned} &2(K(\mathbf{t} + \mathbf{h}) - K(\mathbf{t})) \\ &= 2K(\mathbf{0}) - 2K(\mathbf{t}) - (2K(\mathbf{0}) - 2K(\mathbf{t} + \mathbf{h})) \\ &= \|Y(\mathbf{x}^{(0)} + \mathbf{t}) - Y(\mathbf{x}^{(0)})\|^2 - \|Y(\mathbf{x}^{(0)} + \mathbf{t} + \mathbf{h}) - Y(\mathbf{x}^{(0)})\|^2 \\ &= -\|Y(\mathbf{x}^{(0)} + \mathbf{t} + \mathbf{h}) - Y(\mathbf{x}^{(0)} + \mathbf{t})\|^2 \\ &\quad + 2\mathbb{E} \left[\left(Y(\mathbf{x}^{(0)} + \mathbf{t} + \mathbf{h}) - Y(\mathbf{x}^{(0)} + \mathbf{t}) \right) \left(Y(\mathbf{x}^{(0)} + \mathbf{t}) - Y(\mathbf{x}^{(0)}) \right) \right] \\ &= -\|Y(\mathbf{x}^{(0)} + \mathbf{h}) - Y(\mathbf{x}^{(0)})\|^2 \\ &\quad + 2\mathbb{E} \left[\left(Y(\mathbf{x}^{(0)} + \mathbf{t} + \mathbf{h}) - Y(\mathbf{x}^{(0)} + \mathbf{t}) \right) \left(Y(\mathbf{x}^{(0)} + \mathbf{t}) - Y(\mathbf{x}^{(0)}) \right) \right]. \end{aligned} \quad (1.20)$$

The last equality holds because of Y 's stationarity.

The definition of mean square continuity yields the result. \square

Mean square continuous differentiability of stationary Gaussian processes can also be characterized in terms of the covariance kernel [Bachoc, 2013a, chapter 2].

Proposition 1.25. *A stationary Gaussian process Y is mean square continuously differentiable if, and only if, its covariance kernel K is twice continuously differentiable on $\mathbb{R}_{\mathcal{D}}^r$. If this is the case, for all vectors $\mathbf{u} \in \mathbb{R}^r$, $\partial_{\mathbf{u}}\partial_{\mathbf{u}}K$ is the covariance kernel of the stationary mean square continuous Gaussian process $\partial_{\mathbf{u}}Y$.*

Proof. Assume the stationary process Y is mean square continuously differentiable. We prove that for every couple of vectors $(\mathbf{u}, \mathbf{v}) \in \mathbb{R}^r \times \mathbb{R}^r$, the partial derivative $\partial_{\mathbf{u}}\partial_{\mathbf{v}}K$ exists and is continuous over $\mathbb{R}_{\mathcal{D}}^r$. This is enough to establish that K is twice continuously differentiable on $\mathbb{R}_{\mathcal{D}}^r$.

Equation (1.20) shows that for every $\mathbf{t} \in \mathbb{R}_{\mathcal{D}}^r$

$$\partial_{\mathbf{v}}K(\mathbf{t}) := \lim_{h \rightarrow 0} \frac{K(\mathbf{t} + h\mathbf{v}) - K(\mathbf{t})}{h} = \mathbb{E} \left[\partial_{\mathbf{v}}Y(\mathbf{x}^{(0)}) \left(Y(\mathbf{x}^{(0)} + \mathbf{t}) - Y(\mathbf{x}^{(0)}) \right) \right]. \quad (1.21)$$

So $\partial_{\mathbf{v}}K$ exists and, because Y is mean square continuous on \mathcal{D} , is continuous on $\mathbb{R}_{\mathcal{D}}^r$.

Now, for any $\mathbf{t} \in \mathbb{R}_{\mathcal{D}}^r$, as long as $h > 0$ is small enough, $\mathbf{t} + h\mathbf{u} \in \mathbb{R}_{\mathcal{D}}^r$. We have

$$\begin{aligned} & \partial_{\mathbf{v}}K(\mathbf{t} + h\mathbf{u}) - \partial_{\mathbf{v}}K(\mathbf{t}) \\ &= \mathbb{E} \left[\partial_{\mathbf{v}}Y(\mathbf{x}^{(0)}) \left(Y(\mathbf{x}^{(0)} + \mathbf{t} + h\mathbf{u}) - Y(\mathbf{x}^{(0)}) \right) \right] - \mathbb{E} \left[\partial_{\mathbf{v}}Y(\mathbf{x}^{(0)}) \left(Y(\mathbf{x}^{(0)} + \mathbf{t}) - Y(\mathbf{x}^{(0)}) \right) \right] \\ &= \mathbb{E} \left[\partial_{\mathbf{v}}Y(\mathbf{x}^{(0)}) \left(Y(\mathbf{x}^{(0)} + \mathbf{t} + h\mathbf{u}) - Y(\mathbf{x}^{(0)} + \mathbf{t}) \right) \right]. \end{aligned} \quad (1.22)$$

$$\partial_{\mathbf{u}}\partial_{\mathbf{v}}K(\mathbf{t}) := \lim_{h \rightarrow 0} \frac{\partial_{\mathbf{v}}K(\mathbf{t} + h\mathbf{u}) - \partial_{\mathbf{v}}K(\mathbf{t})}{h} = \mathbb{E} \left[\partial_{\mathbf{u}}Y(\mathbf{x}^{(0)} + \mathbf{t})\partial_{\mathbf{v}}Y(\mathbf{x}^{(0)}) \right]. \quad (1.23)$$

So $\partial_{\mathbf{u}}\partial_{\mathbf{v}}K$ exists and, because $\partial_{\mathbf{v}}Y$ is mean square continuous on \mathcal{D} , is continuous on $\mathbb{R}_{\mathcal{D}}^r$. Notice that Equation (1.23) implies $\partial_{\mathbf{u}}\partial_{\mathbf{u}}K$ is the covariance kernel of the stationary Gaussian process with null mean function $\partial_{\mathbf{u}}Y$.

We now prove the converse. Assume the covariance kernel K of the stationary Gaussian process Y is twice continuously differentiable. Without loss of generality, we assume Y has null mean function. Let $\mathbf{u} \in \mathbb{R}^r$. We prove that $\partial_{\mathbf{u}}Y$ exists and is continuous on \mathcal{D} .

Let $(h_m)_{m \in \mathbb{N}}$ be a sequence of real numbers that converges to 0. First, we show that for all $\mathbf{x}^{(0)} \in \mathcal{D}$ $(h_m^{-1}(Y(\mathbf{x}^{(0)} + h_m\mathbf{u}) - Y(\mathbf{x}^{(0)})))_{m \in \mathbb{N}}$ is a Cauchy sequence in $L^2(\Omega)$.

Let $m, n \in \mathbb{N}$.

$$\begin{aligned} & \left\| \frac{Y(\mathbf{x}^{(0)} + h_m\mathbf{u}) - Y(\mathbf{x}^{(0)})}{h_m} - \frac{Y(\mathbf{x}^{(0)} + h_n\mathbf{u}) - Y(\mathbf{x}^{(0)})}{h_n} \right\|^2 \\ &= \left\| \frac{Y(\mathbf{x}^{(0)} + h_m\mathbf{u}) - Y(\mathbf{x}^{(0)})}{h_m} \right\|^2 + \left\| \frac{Y(\mathbf{x}^{(0)} + h_n\mathbf{u}) - Y(\mathbf{x}^{(0)})}{h_n} \right\|^2 \\ & \quad - \frac{2}{h_m h_n} \mathbb{E} \left[\left(Y(\mathbf{x}^{(0)} + h_m\mathbf{u}) - Y(\mathbf{x}^{(0)}) \right) \left(Y(\mathbf{x}^{(0)} + h_n\mathbf{u}) - Y(\mathbf{x}^{(0)}) \right) \right] \\ &= \frac{2}{h_m^2} (K(\mathbf{0}) - K(h_m\mathbf{u})) + \frac{2}{h_n^2} (K(\mathbf{0}) - K(h_n\mathbf{u})) \\ & \quad - \frac{2}{h_m h_n} [K((h_m - h_n)\mathbf{u}) + K(\mathbf{0}) - K(h_m\mathbf{u}) - K(-h_n\mathbf{u})]. \end{aligned} \quad (1.24)$$

When $m, n \rightarrow \infty$, the third term converges to $2\partial_{\mathbf{u}}\partial_{\mathbf{u}}K(\mathbf{0})$. The first two terms must be examined more closely.

Because K is an even mapping, we have for every $\mathbf{v} \in \mathbb{R}^r$ $K(-\mathbf{v}) = K(\mathbf{v})$. So any partial derivative at $\mathbf{0}$ is necessarily null. Given K is twice continuously differentiable along \mathbf{u} , it admits a Taylor expansion of order 2 in this direction. $K(h_m\mathbf{u}) - K(\mathbf{0}) \underset{m \rightarrow \infty}{\sim} h_m^2 \partial_{\mathbf{u}}\partial_{\mathbf{u}}K(\mathbf{0})/2$, so each of the first two terms converges when $m, n \rightarrow \infty$ to $-\partial_{\mathbf{u}}\partial_{\mathbf{u}}K(\mathbf{0})$. The sum therefore converges to 0 and $(h_m^{-1}(Y(\mathbf{x}^{(0)} + h_m\mathbf{u}) - Y(\mathbf{x}^{(0)})))_{m \in \mathbb{N}}$ is a Cauchy sequence in $L^2(\Omega)$.

$L^2(\Omega)$ is a Banach space so the sequence admits a limit. Let us denote this limit by $\partial_{\mathbf{u}}Y(\mathbf{x}^{(0)})$. It is Gaussian because it is an L^2 limit of Gaussian random variables. Furthermore, the mapping $\partial_{\mathbf{u}}Y : \mathcal{D} \rightarrow L^2(\Omega)$ $\mathbf{x}^{(0)} \mapsto \partial_{\mathbf{u}}Y(\mathbf{x}^{(0)})$ is a Gaussian process because for every

positive integer n and any points $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)} \in \mathcal{D}$, the vector $(\partial_{\mathbf{u}}Y(\mathbf{x}^{(1)}), \dots, \partial_{\mathbf{u}}Y(\mathbf{x}^{(n)}))^{\top}$ is an L^2 limit of Gaussian vectors and is therefore Gaussian. For the same reason, $\partial_{\mathbf{u}}Y$ is stationary. All that remains is to show that it is mean square continuous. This follows from Equation (1.23), which is still valid because it was obtained using nothing more than Y 's stationarity, Y 's mean square continuity and the existence of partial derivatives of Y (but not their continuity). Taking $\mathbf{v} = \mathbf{u}$, it shows that $\partial_{\mathbf{u}}\partial_{\mathbf{u}}K$ is the covariance kernel of $\partial_{\mathbf{u}}Y$. Since it is continuous at $\mathbf{0}$, $\partial_{\mathbf{u}}Y$ is mean square continuous on \mathcal{D} .

□

Corollary 1.26. *For every nonnegative integer n , a stationary Gaussian process Y is of class C^n on \mathcal{D} if, and only if, its covariance kernel K is of class C^{2n} on $\mathbb{R}_{\mathcal{D}}^r$.*

1.4 Spectral representation

Stein [1999] showed how fruitful the spectral representation of Gaussian processes could be in order to understand its properties. In particular, it provides a useful characterization of means square continuity and differentiability.

Bochner [1932]'s theorem is the cornerstone of all spectral results. A few preliminary results, which are by themselves of interest because they provide methods for constructing covariance kernels, are useful for showing it. These results are drawn from Rasmussen and Williams [2006].

In all that follows, $\mathcal{D} = \mathcal{V} = \mathbb{R}^r$ for some positive integer r .

The set of positive definite mappings is stable for several operations listed below.

Proposition 1.27. *The sum of two positive definite mappings $\mathbb{R}^r \rightarrow \mathbb{R}$ is positive definite.*

Proof. This is the translation of the fact that the sum of two positive semi-definite matrices is positive semi-definite. □

Proposition 1.28. *The product of two positive definite mappings $\mathbb{R}^r \rightarrow \mathbb{R}$ is positive definite.*

Proof. Let K and K' be two positive definite mappings $\mathbb{R}^r \rightarrow \mathbb{R}$. Let Y and Y' be independent stationary Gaussian processes with null mean function and covariance kernel K and K' respectively. Let Z be the random process defined by $Z(\mathbf{x}) = Y(\mathbf{x})Y'(\mathbf{x})$ ($\mathbf{x} \in \mathbb{R}^r$).

$$\mathbb{E}[Z(\mathbf{x}^{(1)})Z(\mathbf{x}^{(2)})] = \mathbb{E}[Y(\mathbf{x}^{(1)})Y(\mathbf{x}^{(2)})]\mathbb{E}[Y'(\mathbf{x}^{(1)})Y'(\mathbf{x}^{(2)})] = K(\mathbf{x}^{(1)} - \mathbf{x}^{(2)})K'(\mathbf{x}^{(1)} - \mathbf{x}^{(2)}) \quad (1.25)$$

So the mapping $\mathbf{x} \mapsto K(\mathbf{x})K'(\mathbf{x})$ is a covariance kernel: it is positive definite. □

Proposition 1.29. *Let $K_1 : \mathbb{R}^r \rightarrow \mathbb{R}$ be a continuous function that is integrable with respect to the Lebesgue measure on \mathbb{R}^r and let K_2 be a continuous positive definite mapping $\mathbb{R}^r \rightarrow \mathbb{R}$. Then the “double-convolution” of K_1 and K_2 (defined hereafter) is a positive definite mapping $\mathbb{R}^r \rightarrow \mathbb{R}$.*

For all $\mathbf{x} \in \mathbb{R}^r$,

$$K_1 * K_2 * K_1(\mathbf{x}) = \int_{\mathbb{R}^r} \int_{\mathbb{R}^r} K_1(\mathbf{x} - \mathbf{z})K_2(\mathbf{z} - \mathbf{z}')K_1(\mathbf{z}')d\mathbf{z}d\mathbf{z}' \quad (1.26)$$

Proof. Let Y be a stationary Gaussian process with null mean function and covariance kernel K_2 . First, we need to properly define $Z(\mathbf{x}) := \int_{\mathbb{R}^r} K_1(\mathbf{x} - \mathbf{z})Y(\mathbf{z})d\mathbf{z}$. Since K_2 is continuous, Y is mean square continuous. Since it is stationary, $\|Y(\mathbf{x})\|$ is constant as a function of \mathbf{x} . Moreover, K_1 is integrable. This means that $Z(\mathbf{x}) := \int_{\mathbb{R}^r} K_1(\mathbf{x} - \mathbf{z})Y(\mathbf{z})d\mathbf{z}$ can be defined in the Riemann sense as an L^2 limit of Gaussian random variables and is therefore Gaussian. Similarly, for any positive integer n and any $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$, $(Z(\mathbf{x}^{(1)}), \dots, Z(\mathbf{x}^{(n)}))^\top$ is an L^2 limit of Gaussian vectors and therefore a Gaussian vector, so $Z : \mathbf{x} \mapsto Z(\mathbf{x})$ is a Gaussian process with null mean function.

For all $\mathbf{x}^{(1)}, \mathbf{x}^{(2)} \in \mathbb{R}^r$,

$$\begin{aligned} \mathbb{E}[Z(\mathbf{x}^{(1)})Z(\mathbf{x}^{(2)})] &= \mathbb{E}\left[\int_{\mathbb{R}^r} K_1(\mathbf{x}^{(1)} - \mathbf{z})Y(\mathbf{z})d\mathbf{z} \int_{\mathbb{R}^r} K_1(\mathbf{x}^{(2)} - \mathbf{z}')Y(\mathbf{z}')d\mathbf{z}'\right] \\ &= \int_{\mathbb{R}^r} \int_{\mathbb{R}^r} K_1(\mathbf{x}^{(1)} - \mathbf{z})\mathbb{E}[Y(\mathbf{z})Y(\mathbf{z}')]K_1(\mathbf{x}^{(2)} - \mathbf{z}')d\mathbf{z}d\mathbf{z}' \quad (1.27) \end{aligned}$$

$$\begin{aligned} &= \int_{\mathbb{R}^r} \int_{\mathbb{R}^r} K_1(\mathbf{x}^{(1)} - \mathbf{z})K_2(\mathbf{z} - \mathbf{z}')K_1(\mathbf{x}^{(2)} - \mathbf{z}')d\mathbf{z}d\mathbf{z}' \\ &= \int_{\mathbb{R}^r} \int_{\mathbb{R}^r} K_1(\mathbf{x}^{(1)} - \mathbf{x}^{(2)} - \mathbf{z})K_2(\mathbf{z} - \mathbf{z}')K_1(\mathbf{z}')d\mathbf{z}d\mathbf{z}'. \quad (1.28) \end{aligned}$$

In the above computation, Equation (1.27) holds because the scalar product of the limits is the limit of the scalar products and Equation (1.28) is obtained by the changes of variable $\mathbf{z} \leftarrow \mathbf{z} - \mathbf{x}^{(2)}$ and $\mathbf{z}' \leftarrow \mathbf{z}' - \mathbf{x}^{(2)}$ and because $K_1(-\mathbf{z}') = K_1(\mathbf{z}')$.

Equation (1.28) shows that $K_1 * K_2 * K_1$ is the covariance kernel of a stationary Gaussian process, so it is a positive definite mapping. □

Bochner's theorem is a characterization of positive definite mappings $\mathbb{R}^r \rightarrow \mathbb{R}$ in terms of their spectrum.

Theorem 1.30 (Bochner's theorem). *A mapping $K : \mathbb{R}^r \rightarrow \mathbb{R}$ is continuous and positive definite if, and only if, there exists a finite positive measure μ such that for all $\mathbf{x} \in \mathbb{R}^r$*

$$K(\mathbf{x}) = \int_{\mathbb{R}^r} e^{i\langle \boldsymbol{\omega} | \mathbf{x} \rangle} d\mu(\boldsymbol{\omega}). \quad (1.29)$$

Proof. Let μ be a finite positive measure. For all $\mathbf{x} \in \mathbb{R}^r$, define

$$K(\mathbf{x}) = \int_{\mathbb{R}^r} \exp(i\langle \boldsymbol{\omega} | \mathbf{x} \rangle) d\mu(\boldsymbol{\omega}). \quad (1.30)$$

K is, due to the Dominated Convergence theorem, a continuous mapping $\mathbb{R}^r \rightarrow \mathbb{R}$. We now show it is positive definite.

For every positive integer n and all $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n)^\top \in \mathbb{R}^n$,

$$\begin{aligned} \sum_{i,j \in [1,n]} \xi_i \xi_j K(\mathbf{x}^{(i)} - \mathbf{x}^{(j)}) &= \int_{\mathbb{R}^r} \sum_{i,j \in [1,n]} \left(\xi_i \exp(i\langle \boldsymbol{\omega} | \mathbf{x}^{(i)} \rangle) \overline{\xi_j \exp(i\langle \boldsymbol{\omega} | \mathbf{x}^{(j)} \rangle)} \right) d\mu(\boldsymbol{\omega}) \\ &= \int_{\mathbb{R}^r} \sum_{i=1}^n \left| \xi_i \exp(i\langle \boldsymbol{\omega} | \mathbf{x}^{(i)} \rangle) \right|^2 d\mu(\boldsymbol{\omega}) \\ &\geq 0. \quad (1.31) \end{aligned}$$

To prove the converse result, a step-by-step approach is required. First, we prove the result for positive definite functions that are 1) integrable with respect to the Lebesgue measure on \mathbb{R}^r and 2) whose Fourier transform is also integrable with respect to the Lebesgue measure on \mathbb{R}^r . We then successively relax 2) and 1).

Let K be a positive definite function that satisfies 1) and 2). Then define $\hat{K} : \mathbb{R}^r \rightarrow \mathbb{R}$ by

$$\hat{K}(\boldsymbol{\omega}) = (2\pi)^{-r} \int_{\mathbb{R}^r} K(\boldsymbol{x}) e^{-i\langle \boldsymbol{\omega} | \boldsymbol{x} \rangle} d\boldsymbol{x}. \quad (1.32)$$

By assumption, \hat{K} is integrable. Notice that it is also continuous.

For all $\boldsymbol{x} \in \mathbb{R}^k$, denote by $\phi_{\boldsymbol{x}}$ the mapping $\mathbb{R}^k \rightarrow \mathbb{C}$ defined by $\phi_{\boldsymbol{x}}(\boldsymbol{\omega}) = \exp(i\langle \boldsymbol{\omega} | \boldsymbol{x} \rangle)$.

For any $\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)} \in \mathbb{R}^k$,

$$K(\boldsymbol{x}^{(1)} - \boldsymbol{x}^{(2)}) = \int_{\mathbb{R}^r} \hat{K}(\boldsymbol{\omega}) e^{i\langle \boldsymbol{\omega} | \boldsymbol{x}^{(1)} - \boldsymbol{x}^{(2)} \rangle} d\boldsymbol{\omega} = \int_{\mathbb{R}^r} \phi_{\boldsymbol{x}^{(1)}}(\boldsymbol{\omega}) \overline{\phi_{\boldsymbol{x}^{(2)}}(\boldsymbol{\omega})} \hat{K}(\boldsymbol{\omega}) d\boldsymbol{\omega}. \quad (1.33)$$

So, for any positive integer n , any $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n)^\top \in \mathbb{R}^n$ and any $\boldsymbol{x}^{(1)}, \dots, \boldsymbol{x}^{(n)} \in \mathbb{R}^k$,

$$\begin{aligned} \int_{\mathbb{R}^r} \left| \sum_{k=1}^n \xi_k \phi_{\boldsymbol{x}^{(k)}}(\boldsymbol{\omega}) \right|^2 \hat{K}(\boldsymbol{\omega}) d\boldsymbol{\omega} &= \sum_{k,l \in \llbracket 1, n \rrbracket} \xi_k \xi_l \int_{\mathbb{R}^r} \phi_{\boldsymbol{x}^{(k)}}(\boldsymbol{\omega}) \overline{\phi_{\boldsymbol{x}^{(l)}}(\boldsymbol{\omega})} \hat{K}(\boldsymbol{\omega}) d\boldsymbol{\omega} \\ &= \sum_{k,l \in \llbracket 1, n \rrbracket} \xi_k \xi_l K(\boldsymbol{x}^{(k)} - \boldsymbol{x}^{(l)}) \\ &\geq 0. \end{aligned} \quad (1.34)$$

Equation (1.34) shows that for any element f of the vector space spanned by $\{\phi_{\boldsymbol{x}} : \boldsymbol{x} \in \mathbb{R}^k\}$,

$$\int_{\mathbb{R}^r} |f(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} \in [0, +\infty). \quad (1.35)$$

For all r -tuples \boldsymbol{a} and \boldsymbol{b} such that $\boldsymbol{a} \leq \boldsymbol{b}$ (in the sense that for every $i \in \llbracket 1, r \rrbracket$, $a_i \leq b_i$), define the set $I_{\boldsymbol{a}, \boldsymbol{b}} := [a_1, b_1] \times \dots \times [a_r, b_r]$.

For all r -tuples $\boldsymbol{a}, \boldsymbol{b}, \boldsymbol{a}', \boldsymbol{b}'$ such that $\boldsymbol{a}' \leq \boldsymbol{a} \leq \boldsymbol{b} \leq \boldsymbol{b}'$, define $p_{\boldsymbol{a}, \boldsymbol{b}} : \mathbb{R}^k \rightarrow \mathbb{R}^k$ as the orthogonal projection on $I_{\boldsymbol{a}, \boldsymbol{b}}$ and $p_{\boldsymbol{a}', \boldsymbol{b}'}^c : \mathbb{R}^k \rightarrow \mathbb{R}^k$ as the orthogonal projection on the closure of $\mathbb{R}^k \setminus I_{\boldsymbol{a}', \boldsymbol{b}'}$. Let $f_{\boldsymbol{a}, \boldsymbol{b}, \boldsymbol{a}', \boldsymbol{b}'} : \mathbb{R}^k \rightarrow \mathbb{R}$ be the continuous mapping defined by

$$f_{\boldsymbol{a}, \boldsymbol{b}, \boldsymbol{a}', \boldsymbol{b}'}(\boldsymbol{\omega}) = \frac{\|\boldsymbol{\omega} - p_{\boldsymbol{a}', \boldsymbol{b}'}^c(\boldsymbol{\omega})\|}{\|p_{\boldsymbol{a}, \boldsymbol{b}}(\boldsymbol{\omega}) - p_{\boldsymbol{a}', \boldsymbol{b}'}^c(\boldsymbol{\omega})\|}. \quad (1.36)$$

For all r -tuples $\boldsymbol{a}, \boldsymbol{b}, \boldsymbol{a}', \boldsymbol{b}', \boldsymbol{a}'', \boldsymbol{b}''$ such that $\boldsymbol{a}'' \leq \boldsymbol{a}' \leq \boldsymbol{a} \leq \boldsymbol{b} \leq \boldsymbol{b}' \leq \boldsymbol{b}''$, define the continuous periodic mapping $f_{\boldsymbol{a}, \boldsymbol{b}, \boldsymbol{a}', \boldsymbol{b}', \boldsymbol{a}'', \boldsymbol{b}''} : \mathbb{R}^k \rightarrow \mathbb{R}$ as follows. For every $\boldsymbol{\omega} \in \mathbb{R}^k$, there exists an r -tuple $\boldsymbol{n} = (n_1, \dots, n_r) \in \mathbb{Z}^r$ such that $\boldsymbol{\omega} \in I_{\boldsymbol{a}'' + \boldsymbol{n}, (\boldsymbol{b}'' - \boldsymbol{a}'')}, \boldsymbol{b}' + \boldsymbol{n}, (\boldsymbol{b}' - \boldsymbol{a}'')}$ (here \cdot denotes the dot product). Then

$$f_{\boldsymbol{a}, \boldsymbol{b}, \boldsymbol{a}', \boldsymbol{b}', \boldsymbol{a}'', \boldsymbol{b}''}(\boldsymbol{\omega}) = f_{\boldsymbol{a} + \boldsymbol{n}, (\boldsymbol{b}'' - \boldsymbol{a}''), \boldsymbol{b} + \boldsymbol{n}, (\boldsymbol{b}' - \boldsymbol{a}''), \boldsymbol{a}' + \boldsymbol{n}, (\boldsymbol{b}'' - \boldsymbol{a}''), \boldsymbol{b}' + \boldsymbol{n}, (\boldsymbol{b}' - \boldsymbol{a}'') }(\boldsymbol{\omega}). \quad (1.37)$$

Because $f_{\boldsymbol{a}, \boldsymbol{b}, \boldsymbol{a}', \boldsymbol{b}', \boldsymbol{a}'', \boldsymbol{b}''}$ is a continuous periodic function, there exists a sequence $(f_n)_{n \in \mathbb{N}}$ of functions belonging to the vector space spanned by $\{\phi_{\boldsymbol{x}} : \boldsymbol{x} \in \mathbb{R}^k\}$ that uniformly converges to $f_{\boldsymbol{a}, \boldsymbol{b}, \boldsymbol{a}', \boldsymbol{b}', \boldsymbol{a}'', \boldsymbol{b}''}$. Because the convergence is uniform, there exists a constant $C > 0$ such

that for all $n \in \mathbb{N}$ and all $\boldsymbol{\omega} \in \mathbb{R}^r$, $|f_n(\boldsymbol{\omega})| \leq C$. This fact, along with pointwise convergence makes the Dominated Convergence theorem applicable.

$$\begin{aligned} \int_{\mathbb{R}^r} |f_{\mathbf{a}, \mathbf{b}, \mathbf{a}', \mathbf{b}', \mathbf{a}'', \mathbf{b}''}(\boldsymbol{\omega})|^2 \hat{K}(\boldsymbol{\omega}) d\boldsymbol{\omega} &= \int_{\mathbb{R}^r} \left| \lim_{n \rightarrow \infty} f_n(\boldsymbol{\omega}) \right|^2 \hat{K}(\boldsymbol{\omega}) d\boldsymbol{\omega} \\ &= \lim_{n \rightarrow \infty} \int_{\mathbb{R}^r} |f_n(\boldsymbol{\omega})|^2 \hat{K}(\boldsymbol{\omega}) d\boldsymbol{\omega} \\ &\geq 0. \end{aligned} \quad (1.38)$$

Clearly, for every $\boldsymbol{\omega} \in \mathbb{R}^r$, $|f_{\mathbf{a}, \mathbf{b}, \mathbf{a}', \mathbf{b}', \mathbf{a}'', \mathbf{b}''}(\boldsymbol{\omega})| \leq 1$. Let us make the period $\mathbf{b}'' - \mathbf{a}''$ increase to infinity and apply the Dominated Convergence theorem again.

$$\begin{aligned} \int_{\mathbb{R}^r} |f_{\mathbf{a}, \mathbf{b}, \mathbf{a}', \mathbf{b}'}(\boldsymbol{\omega})|^2 \hat{K}(\boldsymbol{\omega}) d\boldsymbol{\omega} &= \int_{\mathbb{R}^r} \left| \lim_{\substack{\mathbf{a}'' \rightarrow (-\infty)^r \\ \mathbf{b}'' \rightarrow (+\infty)^r}} f_{\mathbf{a}, \mathbf{b}, \mathbf{a}', \mathbf{b}', \mathbf{a}'', \mathbf{b}''}(\boldsymbol{\omega}) \right|^2 \hat{K}(\boldsymbol{\omega}) d\boldsymbol{\omega} \\ &= \lim_{\substack{\mathbf{a}'' \rightarrow (-\infty)^r \\ \mathbf{b}'' \rightarrow (+\infty)^r}} \int_{\mathbb{R}^r} |f_{\mathbf{a}, \mathbf{b}, \mathbf{a}', \mathbf{b}', \mathbf{a}'', \mathbf{b}''}(\boldsymbol{\omega})|^2 \hat{K}(\boldsymbol{\omega}) d\boldsymbol{\omega} \\ &\geq 0. \end{aligned} \quad (1.39)$$

Finally, let us make $\mathbf{a}' \rightarrow \mathbf{a}$ ($\mathbf{a}' \leq \mathbf{a}$) and $\mathbf{b}' \rightarrow \mathbf{b}$ ($\mathbf{b}' \geq \mathbf{b}$) and apply the Dominated Convergence theorem one last time.

$$\begin{aligned} \int_{I_{\mathbf{a}, \mathbf{b}}} \hat{K}(\boldsymbol{\omega}) d\boldsymbol{\omega} &= \int_{\mathbb{R}^r} |\mathbf{1}_{I_{\mathbf{a}, \mathbf{b}}}(\boldsymbol{\omega})|^2 \hat{K}(\boldsymbol{\omega}) d\boldsymbol{\omega} = \int_{\mathbb{R}^r} \left| \lim_{\substack{\mathbf{a}' \rightarrow \mathbf{a} \\ \mathbf{b}' \rightarrow \mathbf{b}}} f_{\mathbf{a}, \mathbf{b}, \mathbf{a}', \mathbf{b}'}(\boldsymbol{\omega}) \right|^2 \hat{K}(\boldsymbol{\omega}) d\boldsymbol{\omega} \\ &= \lim_{\substack{\mathbf{a}' \rightarrow \mathbf{a} \\ \mathbf{b}' \rightarrow \mathbf{b}}} \int_{\mathbb{R}^r} |f_{\mathbf{a}, \mathbf{b}, \mathbf{a}', \mathbf{b}'}(\boldsymbol{\omega})|^2 \hat{K}(\boldsymbol{\omega}) d\boldsymbol{\omega} \\ &\geq 0. \end{aligned} \quad (1.40)$$

Because \hat{K} is continuous, this implies that \hat{K} is a nonnegative function. It is therefore the density with respect to the Lebesgue measure of a positive measure. Moreover, \hat{K} is by assumption integrable with respect to the Lebesgue measure, so it is the density of a finite positive measure.

For all $R > 0$, define the mapping $G_R : \mathbb{R}^r \rightarrow \mathbb{R}$ by $G_R(\mathbf{x}) = (2\pi R^2)^{-r/2} \exp(-\|\mathbf{x}\|^2/(2R^2))$. G_R is continuous and integrable with respect to the Lebesgue measure on \mathbb{R}^r . Its Fourier transform is defined by

$$\hat{G}_R(\boldsymbol{\omega}) = (2\pi)^{-r} \int_{\mathbb{R}^r} G_R(\mathbf{x}) e^{-i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} d\mathbf{x} = (2\pi)^{-r} e^{-\frac{R^2 \|\boldsymbol{\omega}\|^2}{2}}. \quad (1.41)$$

\hat{G}_R is also integrable with respect to the Lebesgue measure on \mathbb{R}^r . It is the density of a positive finite measure, so applying the result shown above, G_R is a positive definite function. It is a ‘‘Squared Exponential’’ kernel.

Squared Exponential kernels are the main tools we will use to relax Assumptions 1) and 2).

First, let us relax Assumption 2). Let K be a continuous positive definite mapping that satisfies Assumption 1): it is integrable with respect to the Lebesgue measure on \mathbb{R}^r .

Then, for all $R > 0$, $G_R * K * G_R$ is also a continuous positive definite mapping. And for all $\boldsymbol{\omega} \in \mathbb{R}^r$,

$$G_R * \widehat{K} * G_R(\boldsymbol{\omega}) = (2\pi)^{-r} \int_{\mathbb{R}^r} G_R * K * G_R(\boldsymbol{x}) e^{-i\langle \boldsymbol{\omega}, \boldsymbol{x} \rangle} d\boldsymbol{x} = \hat{K}(\boldsymbol{\omega}) e^{-R^2 \|\boldsymbol{\omega}\|^2} \quad (1.42)$$

$G_R * \widehat{K} * G_R$ is integrable with respect to the Lebesgue measure on \mathbb{R}^r . Therefore, $G_R * K * G_R$ is a continuous positive definite mapping that satisfies both Assumption 1) and 2). The result established previously can be applied: $G_R * \widehat{K} * G_R$ is the density with respect to the Lebesgue measure on \mathbb{R}^r of a finite positive measure. Equation (1.42) then shows that \hat{K} is a nonnegative mapping, so it is the density of a positive measure.

The Monotone Convergence theorem yields that

$$\lim_{R \rightarrow 0} \int_{\mathbb{R}^r} \hat{K}(\boldsymbol{\omega}) e^{-R^2 \|\boldsymbol{\omega}\|^2} d\boldsymbol{\omega} = \int_{\mathbb{R}^r} \hat{K}(\boldsymbol{\omega}) d\boldsymbol{\omega}. \quad (1.43)$$

For all $R > 0$ $\int_{\mathbb{R}^r} \hat{K}(\boldsymbol{\omega}) e^{-R^2 \|\boldsymbol{\omega}\|^2} d\boldsymbol{\omega} = G_R * K * G_R(\mathbf{0})$. Moreover, because K is continuous and G_R is the probability density of the Normal distribution with null mean and standard deviation R ,

$$\lim_{R \rightarrow 0} G_R * K * G_R(\mathbf{0}) = K(\mathbf{0}). \quad (1.44)$$

Equations (1.43) and (1.44) yield

$$\int_{\mathbb{R}^r} \hat{K}(\boldsymbol{\omega}) d\boldsymbol{\omega} = K(\mathbf{0}), \quad (1.45)$$

so \hat{K} is the density of a finite measure with total mass $K(\mathbf{0})$.

All that remains to do is relax Assumption 1). Let K be a continuous positive definite mapping $\mathbb{R}^r \rightarrow \mathbb{R}$ such that $K(\mathbf{0}) \neq 0$ (because if $K(\mathbf{0}) = 0$, then K is the null function and there is nothing to prove). For all $R > 0$, the mapping $K_R : \boldsymbol{x} \mapsto (2\pi R)^{r/2} K(\boldsymbol{x}) G_R(\boldsymbol{x}) K(\mathbf{0})^{-1}$ is positive definite. Because K is necessarily bounded and G_R is integrable with respect to the Lebesgue measure on \mathbb{R}^r , K_R is integrable too. So, applying the previous result, it is the characteristic function of a probability measure μ_R . When $R \rightarrow +\infty$, K_R converges pointwise to the function $K_0 : \boldsymbol{x} \mapsto K(\boldsymbol{x}) K(\mathbf{0})^{-1}$. This implies that μ_R converges narrowly to a probability measure μ_0 and that K_0 is the characteristic function of μ_0 . Therefore, for all $\boldsymbol{x} \in \mathbb{R}^r$,

$$K(\boldsymbol{x}) = K(\mathbf{0}) \int_{\mathbb{R}^r} e^{i\langle \boldsymbol{\omega}, \boldsymbol{x} \rangle} d\mu_0(\boldsymbol{\omega}). \quad (1.46)$$

□

An other proof can be found in Gihman and Skorohod [1974] page 208.

Sometimes, a more precise version of Bochner's theorem is required. The usual terminology is somewhat misleading. A positive definite kernel leads to positive **semi**-definite covariance matrices. But when does it lead to positive definite covariance matrices?

Proposition 1.31. *Let μ be a positive measure on \mathbb{R}^r with finite non-null total mass that is absolutely continuous with respect to the Lebesgue measure. Then the mapping $K : \mathbb{R}^r \rightarrow \mathbb{R}$ defined by*

$$K(\mathbf{x}) = \int_{\mathbb{R}^r} e^{i\langle \boldsymbol{\omega} | \mathbf{x} \rangle} d\mu(\boldsymbol{\omega}) \quad (1.47)$$

is positive definite. Moreover, for any $\boldsymbol{\xi} \in \mathbb{R}^n \setminus \{\mathbf{0}_n\}$,

$$\sum_{k,l \in [1,n]} \xi_k \xi_l K(\mathbf{x}^{(k)} - \mathbf{x}^{(l)}) > 0. \quad (1.48)$$

Proof. The first part results from Bochner's theorem. Let us show the second.

$$\begin{aligned} \sum_{k,l \in [1,n]} \xi_k \xi_l K(\mathbf{x}^{(k)} - \mathbf{x}^{(l)}) &= \sum_{k,l \in [1,n]} \xi_k \xi_l \int_{\mathbb{R}^r} e^{i\langle \boldsymbol{\omega} | \mathbf{x}^{(k)} - \mathbf{x}^{(l)} \rangle} d\mu(\boldsymbol{\omega}) \\ &= \int_{\mathbb{R}^r} \left| \sum_{k=1}^n \xi_k e^{i\langle \boldsymbol{\omega} | \mathbf{x}^{(k)} \rangle} \right|^2 d\mu(\boldsymbol{\omega}). \end{aligned} \quad (1.49)$$

Given $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$ are all distinct, for almost all unitary vector \mathbf{u} in the sense of the Lebesgue measure on the unit sphere S^{r-1} , the real numbers $\langle \mathbf{u} | \mathbf{x}^{(1)} \rangle, \dots, \langle \mathbf{u} | \mathbf{x}^{(n)} \rangle$ are distinct. Indeed, if two of these numbers, say $\langle \mathbf{u} | \mathbf{x}^{(1)} \rangle$ and $\langle \mathbf{u} | \mathbf{x}^{(2)} \rangle$, were equal, then \mathbf{u} would be orthogonal to $\mathbf{x}^{(1)} - \mathbf{x}^{(2)}$. But the set of all vectors of \mathbb{R}^r orthogonal to $\mathbf{x}^{(1)} - \mathbf{x}^{(2)}$ is a hyperplane. So there exists a finite number of hyperplanes of \mathbb{R}^r such that, if \mathbf{u} does not belong to any of them, the real numbers $\langle \mathbf{u} | \mathbf{x}^{(1)} \rangle, \dots, \langle \mathbf{u} | \mathbf{x}^{(n)} \rangle$ are distinct.

Now, notice that the mapping $\mathbb{C} \rightarrow \mathbb{C}; z \mapsto \sum_{k=1}^n \xi_k e^{iz\langle \mathbf{u} | \mathbf{x}^{(k)} \rangle}$ is holomorphic. So either it is the null function or all its zeros are isolated. Given it clearly is not the null function, its zeros are isolated. So the set of all zeros that belong to \mathbb{R} is countable and therefore of null Lebesgue measure.

Let $f : \mathbb{R}^r \rightarrow \{0, 1\}$ be the measurable mapping such that $f(\boldsymbol{\omega}) = 1$ if $\sum_{k=1}^n \xi_k e^{i\langle \boldsymbol{\omega} | \mathbf{x}^{(k)} \rangle} = 0$ and $f(\boldsymbol{\omega}) = 0$ if not.

$$\int_{\mathbb{R}^r} f(\boldsymbol{\omega}) d\boldsymbol{\omega} = \frac{2\pi^{\frac{r}{2}}}{\Gamma(\frac{r}{2})} \int_{S^{r-1}} \int_{(0,+\infty)} f(t\mathbf{u}) t^{r-1} dt d\mathbf{u} = \frac{2\pi^{\frac{r}{2}}}{\Gamma(\frac{r}{2})} \int_{S^{r-1}} 0 d\mathbf{u} = 0. \quad (1.50)$$

Therefore the mapping $\mathbb{R}^r \rightarrow \mathbb{C}; \boldsymbol{\omega} \mapsto \sum_{k=1}^n \xi_k e^{i\langle \boldsymbol{\omega} | \mathbf{x}^{(k)} \rangle}$ takes null values on a Borel set that is negligible with respect to the Lebesgue measure. This set is therefore also negligible with respect to μ , which yields the conclusion. \square

The following results give the link between the existence of moments of the spectral measure and the smoothness of any associated Gaussian process [Stein, 1999, chapter 2].

Proposition 1.32. *For every nonnegative integer n , a positive definite mapping $\mathbb{R}^r \rightarrow \mathbb{R}$ is of class C^n if, and only if, its Fourier transform μ verifies*

$$\int_{\mathbb{R}^r} \|\boldsymbol{\omega}\|^n d\mu(\boldsymbol{\omega}) < +\infty. \quad (1.51)$$

Corollary 1.33. *For every nonnegative integer n , a stationary Gaussian process $Y : \mathcal{D} \subset \mathbb{R}^k \rightarrow L^2(\Omega)$ whose covariance kernel K is positive definite on \mathbb{R}^k is of class C^n on \mathcal{D} if, and only if, K 's Fourier transform μ verifies*

$$\int_{\mathbb{R}^r} \|\boldsymbol{\omega}\|^{2n} d\mu(\boldsymbol{\omega}) < +\infty. \quad (1.52)$$

To illustrate this, we provide in Figure 1.2 Kriging examples for differently smooth processes.

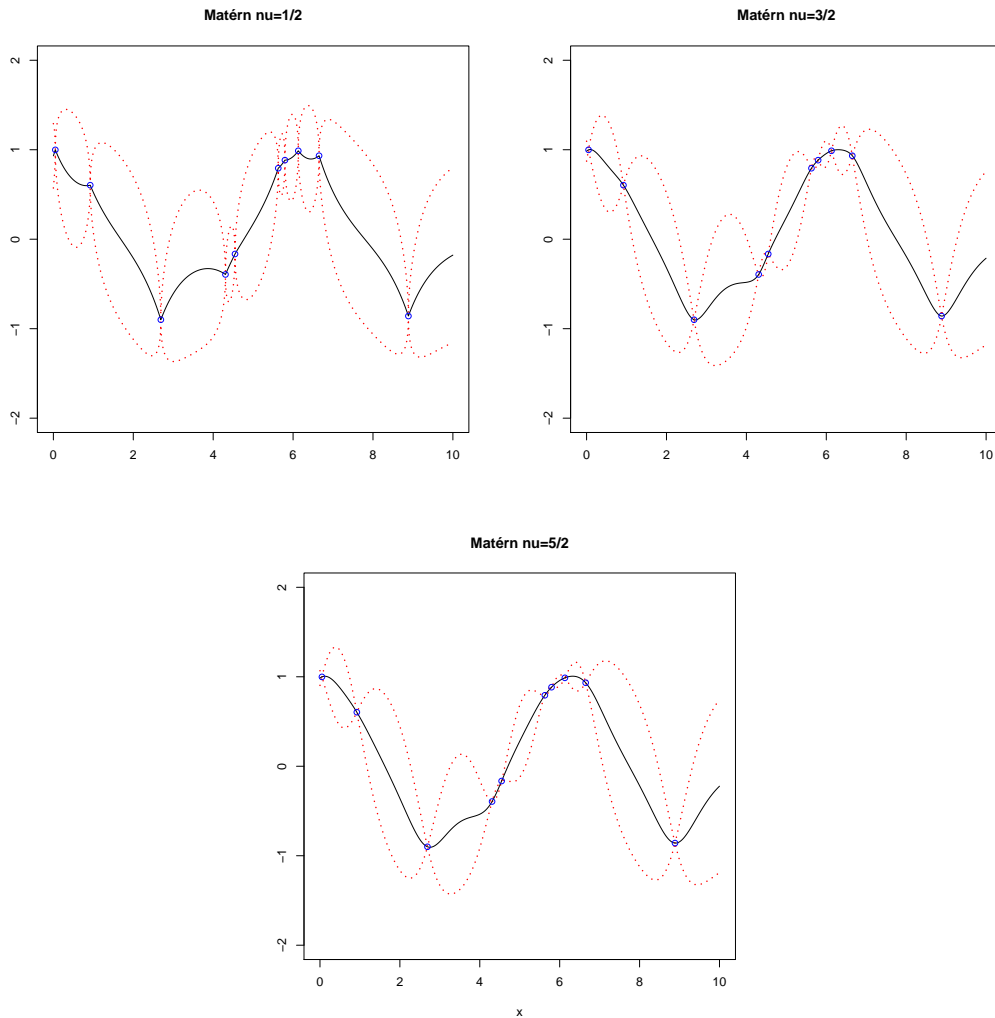


Figure 1.2 – *Kriging with differently smooth Gaussian processes. Top left: mean square continuous. Top right: mean square continuously differentiable. Bottom: twice mean square continuously differentiable. Solid lines represent conditional means and dotted lines add or subtract conditional standard deviations.*

1.5 Examples of correlation kernels

In this thesis, four families of correlation kernels are discussed (cf. Table 1.1). All corresponding spectral probability measures are absolutely continuous with respect to the Lebesgue measure.

Kernel	$K_\theta(x)$	parameter range
Spherical ($r = 1, 2, 3$)	$\left(1 - \frac{3}{2} \left(\frac{ x }{\theta}\right) + \frac{1}{2} \left(\frac{ x }{\theta}\right)^3\right) \mathbf{1}_{\{ x \leq \theta\}}$	\emptyset
Power Exponential	$\exp\left\{-\left(\frac{ x }{\theta}\right)^q\right\}$	$q \in (0, 2]$
Rational Quadratic	$\left(1 + \left(\frac{ x }{\theta}\right)^2\right)^{-\nu}$	$\nu \in (0, +\infty)$
Matérn	$\Gamma(\nu)^{-1} 2^{1-\nu} \left(2\sqrt{\nu} \frac{ x }{\theta}\right)^\nu \mathcal{K}_\nu\left(2\sqrt{\nu} \frac{ x }{\theta}\right)$	$\nu \in (0, +\infty)$

Table 1.1 – Formulas for several correlation kernel families. The Squared Exponential kernel is the Power Exponential kernel with $q = 2$. \mathcal{K}_ν is the modified Bessel function of second kind with parameter ν [Abramowitz and Stegun, 1964](9.6.).

The spherical kernel is only used in geostatistics ([Journel and Huijbregts, 1978], p.116, [Isaaks and Srivastava, 1989], p.374, [Bras and Rodriguez-Iturbe, 1985], p.418, [Christakos, 1992], p.71, [Wackernagel, 1995], p.42, [Kitadinis, 1997], p.56, [Goovaerts, 1997], p.88) mainly because it is limited to 1, 2 and 3-dimensional settings: Bochner’s theorem can be used to show that its Fourier transform in 4 or higher-dimensional settings is no positive measure. Its main advantage is its simplicity. It is continuous, so Gaussian processes using it are mean square continuous. It is not differentiable at 0, so corresponding Gaussian processes are not mean square continuously differentiable. Moreover, the values of the Gaussian process at points with distance greater or equal to θ are independent. This is a peculiar assumption that is generally hard to justify, but can lead to sparsity in correlation matrices.

Spectral density depends on the dimension. In the following, J_0 and J_1 are respectively the Bessel function of first kind with parameter 0 and 1. H_0 and H_1 are respectively the Struve functions with parameter 0 and 1.

$$\hat{K}(\boldsymbol{\omega}) = \begin{cases} \theta \frac{3(\|\theta\boldsymbol{\omega}\|^2 - 2\|\theta\boldsymbol{\omega}\| \frac{\sin(\|\theta\boldsymbol{\omega}\|)}{2\pi\theta\|\theta\boldsymbol{\omega}\|^4} - 2\cos(\|\theta\boldsymbol{\omega}\|) + 2)}{2\pi\theta\|\theta\boldsymbol{\omega}\|^4} & \text{for } r = 1. \\ \theta^2 \left[\frac{H_0(\|\theta\boldsymbol{\omega}\|)J_1(\|\theta\boldsymbol{\omega}\|)}{4} + \frac{(2 - \pi H_1(\|\theta\boldsymbol{\omega}\|))J_0(\|\theta\boldsymbol{\omega}\|)}{4\pi} - \frac{J_1(\|\theta\boldsymbol{\omega}\|)}{\|\theta\boldsymbol{\omega}\|} - \frac{J_2(\|\theta\boldsymbol{\omega}\|)}{\|\theta\boldsymbol{\omega}\|^2} \right] & \text{for } r = 2. \\ \theta^3 \frac{\|\theta\boldsymbol{\omega}\|^3 - 3\sin(\|\theta\boldsymbol{\omega}\|) + 3\|\theta\boldsymbol{\omega}\| \cos(\|\theta\boldsymbol{\omega}\|)}{\pi^2\|\theta\boldsymbol{\omega}\|^5} & \text{for } r = 3. \end{cases}$$

Power Exponential kernels, and especially the Exponential ($q = 1$) and Squared Exponential ($q = 2$) kernels are probably the most widely used family. They are as simple as the Spherical kernel and have the advantage of extending to dimensions greater than 3. They provide some flexibility with respect the smoothness of the associated Gaussian processes, although it is of the all-or-nothing kind. For $q < 2$, Gaussian processes are mean square continuous but not mean square differentiable. For $q = 2$, Gaussian processes are infinitely mean square continuously differentiable. In fact, Stein [1999] points out that in a 1-dimensional setting, letting $Y^{(i)}(0)$ be the i -th derivative of the Gaussian process at 0, $\sum_{i=0}^n Y^{(i)}(0)t^i/i!$ converges to $Y(t)$ in L^2 . So observing Y on any neighborhood of 0 (and therefore of any point) is enough to almost surely know the value of Y at any other point! This behavior is usually considered physically unrealistic.

The Squared Exponential kernel has an other remarkable feature: it is left invariant (up to a multiplicative constant) by the Fourier transform:

$$\hat{K}(\boldsymbol{\omega}) = (2\pi)^{-r} (\pi\theta^2)^{r/2} \exp(-\theta^2\|\boldsymbol{\omega}\|^2/4). \quad (1.53)$$

Rational Quadratic kernels provide no flexibility with respect to smoothness: all lead to infinitely mean square continuously differentiable Gaussian processes. They share the Squared Exponential kernel's defect, but on a smaller scale: $\sum_{i=0}^n Y^{(i)}(0)t^i/i!$ converges to $Y(t)$ in L^2 only if $t < \theta$.

Rational Quadratic kernels have the following spectral density, with \mathcal{K}_ν being the modified Bessel function of second kind with parameter ν :

$$\hat{K}(\boldsymbol{\omega}) = \theta^r \frac{(2\pi)^{-r/2} \|\boldsymbol{\omega}\|^{\nu-r/2} \mathcal{K}_{\nu-r/2}(\|\boldsymbol{\omega}\|)}{2^{\nu-1} \Gamma(\nu)}. \quad (1.54)$$

Matérn kernels are recommended by Stein [1999] because they offer great flexibility regarding smoothness of the Gaussian process. The density with respect to the Lebesgue measure of the associated spectral probability measure is proportional to $(4\nu + \|\boldsymbol{\omega}\|^2/\theta^2)^{-\frac{r}{2}-\nu}$ [Rasmussen and Williams, 2006]. In fact, this kernel is specifically constructed in order to have a nice spectral density. The study of this measure is straightforward and leads to the conclusion that the associated Gaussian process is $\lfloor \nu \rfloor$ times mean square continuously differentiable if ν is not an integer and $\nu - 1$ continuously differentiable if it is. Moreover, it is mean square continuous for any $\nu > 0$. The Matérn family admits comfortable expressions when ν is a half-integer ($\nu = 1/2, 3/2, 5/2, \dots$). Notably, it includes the Exponential kernel ($\nu = 1/2$) and admits the Squared Exponential kernel as its limit when $\nu \rightarrow +\infty$. Notice that Matérn and Rational Quadratic kernels are the Fourier transform of one another.

Parameter ν	Analytical expression
$\frac{1}{2}$	$\sigma^2 e^{-\sqrt{2} \frac{ x }{\ell}}$
$\frac{3}{2}$	$\sigma^2 \left(1 + \sqrt{6} \frac{ x }{\ell}\right) e^{-\sqrt{6} \frac{ x }{\ell}}$
$\frac{5}{2}$	$\sigma^2 \left(1 + \sqrt{10} \frac{ x }{\ell} + \frac{10}{3} \frac{ x ^2}{\ell^2}\right) e^{-\sqrt{10} \frac{ x }{\ell}}$
∞	$\sigma^2 e^{-\frac{ x ^2}{\ell^2}}$

Figure 1.3 illustrates several kernels from all four families.

Estimating kernel parameters is often no simple task, and is one of the motivations of this thesis. Here, the matter is tackled from an Objective Bayesian standpoint. For a frequentist perspective, see Bachoc [2013b].

1.6 Current Kriging-related research

Although the aim of this chapter was only to provide an understanding of the Kriging tools used in this thesis, let us briefly mention several recent developments in the field. See also Bachoc et al. [2017a].

Apart from geostatistical applications, one of the main uses for Kriging methods is emulation of computer codes. Kriging is therefore used for both calibration and prediction purposes [Bachoc et al., 2014]. Gaussian processes can be used to perform uncertainty quantification when optimizing under constraints [Binois et al., 2015]. The Bayesian framework allows additional flexibility by considering several different covariance kernels at once [Pronzato and Rendas, 2017] – the present thesis makes such an approach more systematic.

An important issue for all spatial statistics is the design of experiments. One can for example wish for a “space-filling” design set. Although the concept is simple, practical implementation

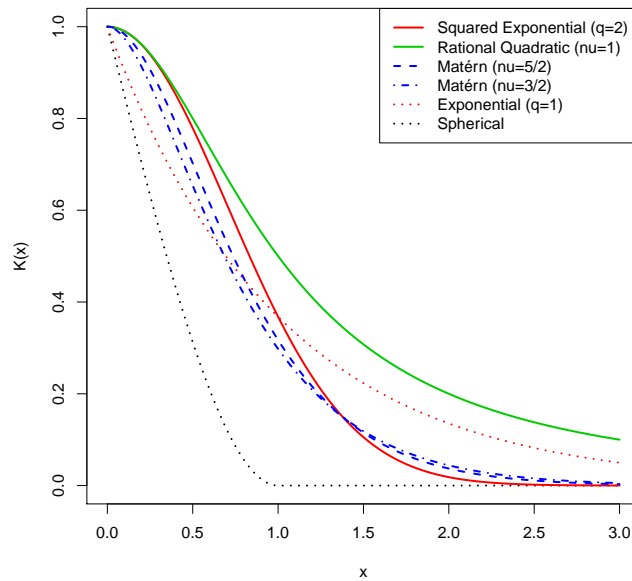


Figure 1.3 – Correlation kernels with parameter $\theta = 1$.

in higher-dimensional spaces is a difficult problem [Damblin et al., 2013, Pronzato, 2017]. Other approaches rely on the statistical properties of Gaussian processes in order to minimize global uncertainty [Chevalier et al., 2014, Bect et al., 2016]. In some applications, a spherical symmetry may be observed and Kriging methods would need to take this information into account [Padonou and Roustant, 2016].

Because of its versatility, Kriging can be used in varied industrial problems. The one which motivated this thesis is the computation of Functional Risk Curves. Although by no means the only possible solution, the Kriging framework is particularly well suited because it provides ready-to-use uncertainty quantifying objects like prediction intervals [Iooss and Le Gratiet, 2017, Le Gratiet et al., 2017].

Despite its original geostatistical purpose of quantifying quantities in 3-dimensional spaces, it is part of the solution to ever more complex problems. Co-Kriging, i.e. Kriging with multidimensional output, has been successfully used [Le Gratiet and Garnier, 2014] to emulate hierarchical multi-fidelity codes [Le Gratiet, 2013]. In a different context, methods have been developed to deal with nested codes, where the output of one is the input of the next [Perrin et al., 2017]. And when the number of data is large, small Kriging models can be aggregated [Rulli ere et al., 2018].

Qualitative limits of Kriging are also being pushed. Roustant et al. [2018] consider Kriging-based methods for categorical inputs, and Bachoc et al. [2017b] propose a theory for Kriging of distributional rather than numerical inputs.

Chapter 2

Reference Prior Theory

Abstract

Reference analysis is presented with inferential and predictive goals in mind (hypothesis testing is not considered). Formal definitions of reference priors are given for single-parameter and multi-parameter models. Reference priors are linked to both Laplace's insufficient reason argument and Jeffrey's rule for prior elicitation. Fundamental properties of reference priors – independence of sample size, compatibility with sufficient statistics, invariance under reparametrization – are highlighted. Although primarily a review on reference analysis in the vein of Bernardo [2005] and Berger et al. [2009], the chapter contributes a uniqueness result under some assumptions for single-parameter reference priors. Examples of single-parameter and multi-parameter models and associated reference priors are given.

Résumé

L'analyse bayésienne fondée sur les priors de référence est présentée dans une optique inférentielle et prédictive (le test d'hypothèses n'est pas traité). Les priors de référence pour modèles à un ou plusieurs paramètres sont formellement définis. La théorie est reliée à la fois au principe de raison insuffisante de Laplace et à la règle de Jeffreys pour éliciter des priors. Les propriétés fondamentales des priors de référence – indépendance vis-à-vis de la taille d'échantillon, compatibilité avec les statistiques exhaustives, invariance par reparamétrisation – sont mises en relief. Quoique se voulant avant tout un résumé de la théorie dans les pas de Bernardo [2005] et Berger et al. [2009], ce chapitre apporte une contribution sous la forme d'un résultat d'unicité sous conditions du prior de référence pour les modèles à un seul paramètre. Des exemples de modèles à un ou plusieurs paramètres sont fournis avec les priors de référence qui y sont associés.

2.1 Introduction

This chapter aims to provide the basics necessary to gain a working understanding of reference prior theory. It draws mainly on Bernardo [2005] and Berger et al. [2009]. Reference prior theory is an attempt to create “objective” priors. Such priors should be sensible defaults in cases where no prior information is available.

The original noninformative prior was the uniform prior. For finite states problems, it derives from the insufficient reason principle ascribed to Laplace [Gelman et al., 2013]. Its extension

to continuous state spaces is not defensible, however, because it leads to widely different posterior distributions for different parametrizations. Subsequent approaches to the problems have relied on invariance arguments (see Hartigan [1964], Jaynes [1968], Dawid [1983]), resulting in left and right Haar measures [Nachbin, 1965]. Jaynes based the notion of non-informativity on Shannon entropy [Jaynes, 1982]. Despite great successes in many problems like spectral analysis and image processing, this approach was shown to lead to paradoxical results [Seidenfeld, 1986]. There have also been efforts to find priors that make Bayesian credible intervals fit frequentist confidence intervals, starting with Lindley [1958] and Welch and Peers [1963]. For a review of efforts to find noninformative priors, see Kass and Wasserman [1996] or [Ghosh, 2011]. The most fruitful approach however was the Jeffreys-rule prior [Jeffreys, 1961], and many methods end up justifying its use in several cases. As we will see, the reference prior method is one of them.

Because of the vague nature of the notion of noninformative prior – each method is a formalization of this idea – it is the subject of an ongoing controversy. We do not discuss it in this chapter, but the interested reader may for example read Fienberg [2006] and Berger [2006], or more recently Seaman III et al. [2012] and Kamary and Robert [2014].

Apart from inference and prediction, hypothesis testing is yet another statistical problem where noninformative priors are desired. Although no such problem is considered in this thesis, it is a very active research field. See Bernardo [2011] for solutions based on reference prior theory, and Kamary et al. [2014] for a general modern perspective.

This chapter details the intuition behind the formal definition of the reference prior in order to explain both its strengths and shortcomings. The formal definition is not given straight away, because it is quite complex and is of little help if the ideas behind it are not well understood. Instead, it is presented as the end result of a development process where “naive” conceptions are put forward, criticized, and then corrected.

Even though intuition is emphasized, formal proofs are provided for all results except Theorem 2.16, for which heuristic arguments are provided. The reader is referred to Clarke and Barron [1994] for formal arguments in this instance.

2.2 Basic idea

Our main tool for defining the reference prior is the Kullback-Leibler divergence.

Definition 2.1 (Kullback-Leibler divergence). *Let P and Q be two probability measures absolutely continuous with respect to a measure μ . Let p and q be their respective Radon-Nikodym derivatives with respect to μ . The Kullback-Leibler divergence $D(P||Q)$ of Q with respect to P is defined by*

$$D(P||Q) = \int p \log \left(\frac{p}{q} \right) d\mu. \quad (2.1)$$

The Kullback-Leibler divergence has several useful properties [Lindley, 1956]. Among them:

- It does not depend on the dominating measure μ .
- Jensen’s inequality can be used to show that it is nonnegative.
- It is equal to 0 if, and only if, $Q = P$.

$$- \int p \left| \log \left(\frac{p}{q} \right) \right| d\mu \leq D(P||Q) + 1.$$

The Kullback-Leibler divergence is usually interpreted as a way to quantify how different Q is from P . The smaller $D(P||Q)$ is, the better Q approximates P . In the discrete case, it is axiomatically justified in the context of Shannon's information theory [Shannon, 1948, Lee, 1964]. Good [1966] gives a probabilistic interpretation of information. Bernardo [1979b] discusses it from a decision-theoretic standpoint.

The Kullback-Leibler divergence is not symmetric, so it does not qualify as a distance between probability distributions.

This interpretation can be used to determine what an objective prior distribution should be. Let \mathcal{Y} be the observation space and endow it with a σ -algebra \mathbb{Y} . Let Θ be the parameter space. A model is a collection of probability distributions $(P_\theta)_{\theta \in \Theta}$. In a Bayesian setting, it is coupled with a prior distribution Π on the parameter space Θ . For ease of use, we assume that Θ is a metric space (its distance is denoted by *dist*) endowed with its Borel σ -algebra. Every prior distribution Π is a probability distribution on this measurable space. The predictive distribution based on prior knowledge [Robert, 2007] is

$$P_\Pi := \int_{\Theta} P_\theta d\Pi(\theta). \quad (2.2)$$

This is to say that for any measurable subset T of \mathcal{Y} ,

$$P_\Pi(T) = \int_{\Theta} P_\theta(T) d\Pi(\theta). \quad (2.3)$$

Naturally, this definition assumes that for any measurable set $A \in \mathbb{Y}$, the mapping $\Theta \rightarrow [0, 1]; \theta \mapsto P_\theta(A)$ is measurable. In other words, the mapping $\Theta \times \mathbb{Y}; (\theta, A) \mapsto P_\theta(A)$ is assumed to be a Markov kernel. This assumption is technical and has no practical relevance. Typically, $P_\theta(A)$ is continuous as a function of θ , or else Θ is a countable set, or some other reason makes this requirement hold.

If an observation $y \in \mathcal{Y}$ from the model $(P_\theta)_{\theta \in \Theta}$ (with unknown θ) has been made, let $\Pi(\cdot|y)$ be the posterior distribution resulting from the prior Π .

The posterior distribution [Robert, 2007] is defined as a Markov kernel $\mathcal{Y} \times \mathcal{B}(\Theta) : (y, B) \mapsto \Pi(B|y)$ such that for any measurable set $A \subset \mathcal{Y}$ and any measurable set $B \subset \Theta$

$$\int_A \Pi(B|y) dP_\Pi(y) = \int_B P_\theta(A) d\Pi(\theta). \quad (2.4)$$

Remark. If the model $(P_\theta)_{\theta \in \Theta}$ is dominated by some measure $\mu_{\mathcal{Y}}$ (i.e. if for any $\theta \in \Theta$, P_θ is absolutely continuous with respect to $\mu_{\mathcal{Y}}$ – in other words, if the model admits a likelihood function), then let p_θ be the density of P_θ and p_Π be the density of P_Π (for $\mu_{\mathcal{Y}}$ -almost any y , $p_\Pi(y) = \int_{\Theta} p_\theta(y) d\Pi(\theta)$). Then, for $\mu_{\mathcal{Y}}$ -almost any y , the posterior distribution $\Pi(\cdot|y)$ is defined as having density $\theta \mapsto p_\theta(y)/p_\Pi(y)$ with respect to the prior Π .

Intuitively, a prior Π is noninformative if it plays a small role in the formation the posterior $\Pi(\cdot|y)$. This way, the amount of information conveyed by the data will be maximized. In practice, we use the Kullback-Leibler divergence $D(\Pi(\cdot|y) || \Pi)$, which says how badly Π approximates $\Pi(\cdot|y)$ to measure this [Berger et al., 2009]. Given that y is actually unknown at the time the prior is defined, $D(\Pi(\cdot|y) || \Pi)$ should be maximized on average over $y \in \mathcal{Y}$. We therefore want to maximize over all priors Π the quantity $\int_{\mathcal{Y}} dP_\Pi(y) D(\Pi(\cdot|y) || \Pi)$.

If the model is dominated (cf. previous remark), then

$$\begin{aligned} \int_{\mathcal{Y}} dP_{\Pi}(y)D(\Pi(\cdot|y)||\Pi) &= \int_{\mathcal{Y}} dP_{\Pi}(y) \int_{\Theta} d\Pi(\theta|y) \log \left(\frac{d\Pi(\cdot|y)}{d\Pi}(\theta) \right) \\ &= \int_{\mathcal{Y}} dP_{\Pi}(y) \int_{\Theta} d\Pi(\theta|y) \log \frac{p_{\theta}(y)}{p_{\Pi}(y)}. \end{aligned} \quad (2.5)$$

If this quantity is finite, then the Fubini-Lebesgue theorem is applicable, because

$$\int_{\mathcal{Y}} dP_{\Pi}(y) \int_{\Theta} d\Pi(\theta|y) \left| \log \frac{p_{\theta}(y)}{p_{\Pi}(y)} \right| \leq \int_{\mathcal{Y}} dP_{\Pi}(y) \int_{\Theta} d\Pi(\theta|y) \log \frac{p_{\theta}(y)}{p_{\Pi}(y)} + 1. \quad (2.6)$$

Therefore, if $\int_{\mathcal{Y}} dP_{\Pi}(y)D(\Pi(\cdot|y)||\Pi) < +\infty$, then the two integral signs can be switched.

$$\int_{\mathcal{Y}} dP_{\Pi}(y)D(\Pi(\cdot|y)||\Pi) = \int_{\Theta} d\Pi(\theta) \int_{\mathcal{Y}} p_{\theta}(y) \log \frac{p_{\theta}(y)}{p_{\Pi}(y)} d\mu_{\mathcal{Y}}(y) = \int_{\Theta} d\Pi(\theta)D(P_{\theta}||P_{\Pi}). \quad (2.7)$$

The same reasoning could be applied in reverse. So either both $\int_{\mathcal{Y}} dP_{\Pi}(y)D(\Pi(\cdot|y)||\Pi)$ and $\int_{\Theta} d\Pi(\theta)D(P_{\theta}||P_{\Pi})$ are finite, in which case they are equal, or both are infinite. In all cases, they are equal.

This identity, which holds for all dominated models $(P_{\theta})_{\theta \in \Theta}$, leads to a new intuition. A prior Π is noninformative if there is a great difference between knowing it and knowing the true value θ , that is, if P_{Π} is generally a poor substitute to P_{θ} . So $D(P_{\theta}||P_{\Pi})$ should be large on average over Θ .

Of course, “on average over Θ ” means in practice “on average over the prior distribution”. Averaging over the prior whose noninformativity is being evaluated makes the noninformativity criterion somewhat self-referential. One would perhaps wish to average over some other measure, but which one? Note that the original idea of maximizing $\int_{\mathcal{Y}} dP_{\Pi}(y)D(\Pi(\cdot|y)||\Pi)$ runs into the same problem. In that formulation, the average is being made over the predictive distribution, but this distribution still depends on the prior! The formulation (2.7) can however be defended from a game-theoretic standpoint [Clarke and Barron, 1994]. Imagine a game between two players, Nature and the Statistician. Nature first picks the parameter θ randomly using a prior distribution Π . Then, using the picked θ , Nature draws y randomly using the distribution P_{θ} . Then, knowing only the prior distribution Π , the Statistician attempts to guess y . For Nature, the optimal play is to maximize (2.7), because on average, this is what makes the Statistician’s job most difficult.

Example 1. Let $\Theta = \{1, 2\}$, P_1 be the uniform distribution on $[0, 1]$ and P_2 be the uniform distribution on $[0, 2]$. Let Π be a prior distribution and let $x := \Pi(1)$ (so that $\Pi(2) = 1 - x$). Then

$$\int_{\Theta} D(P_{\theta}||P_{\Pi})d\Pi(\theta) = x \log \left(\frac{2}{1+x} \right) + \frac{1-x}{2} \log \left(\frac{1}{1-x^2} \right), \quad (2.8)$$

which reaches its maximum at $x = 3/5$.

A problem with this definition is that P_{θ} can be changed simply by taking several independent observations instead of just one. Such a change cannot be said to alter the model, so one would expect the optimum not to change. Unfortunately, it does.

For any distribution Q , let $Q^{\otimes n}$ ($n \in \mathbb{N}$) be the distribution of a sample of n independent draws from Q . And let $P_{\Pi,n}$ be the distribution on \mathcal{Y}^n (with σ -algebra $\mathbb{Y}^{\otimes n}$) such that for all $T_1, \dots, T_n \in \mathcal{Y}$,

$$P_{\Pi,n}(\times_{i=1}^n T_i) = \int_{\Theta} \prod_{i=1}^n P_{\theta}(T_i) d\Pi(\theta). \quad (2.9)$$

Example 2 (Previous example–continued).

$$\int_{\Theta} D(P_{\theta}^{\otimes 2} \| P_{\Pi,2}) d\Pi(\theta) = x \log \left(\frac{4}{1+3x} \right) + (1-x) \left[\frac{1}{4} \log \left(\frac{1}{1+3x} \right) + \frac{3}{4} \log \left(\frac{1}{1-x} \right) \right], \quad (2.10)$$

which reaches its maximum at $x = 1/283(247 - 128 \cdot 2^{1/3} + 48 \cdot 2^{2/3}) \approx 0.57 < 3/5$.

A solution to this problem is to consider the asymptotic case ($n \rightarrow +\infty$).

Let $\mathcal{P}(\Theta)$ be the set all proper prior distributions supported on Θ , that is all probability distributions supported on Θ for which a posterior distribution is well defined. In the case of dominated models, all probability distributions supported on Θ fit this requirement.

We give a first definition of a reference prior, that is a prior which maximizes a noninformativity criterion and can therefore be used as a default prior when little prior information is available.

Definition 2.2 (Naive definition of a reference prior). *A reference prior Π^* is an element of $\mathcal{P}(\Theta)$ such that*

$$\lim_{n \rightarrow +\infty} \int_{\Theta} D(P_{\theta}^{\otimes n} \| P_{\Pi^*,n}) d\Pi^*(\theta) = \lim_{n \rightarrow +\infty} \sup_{\Pi \in \mathcal{P}(\Theta)} \int_{\Theta} D(P_{\theta}^{\otimes n} \| P_{\Pi,n}) d\Pi(\theta) < +\infty. \quad (2.11)$$

Remark. The definition assumes that both limits exist and are finite. As will be seen, this assumption is too strong for practical purposes [Berger et al., 2009]. Nevertheless, there are cases where this definition works [Bernardo, 2005].

Proposition 2.3. *If Θ is a finite set and its topology is the set of all its parts, and if for all $\theta_1, \theta_2 \in \Theta$, $P_{\theta_1} \neq P_{\theta_2}$, then the reference prior is given by*

$$\forall \theta \in \Theta, \quad \Pi^*(\theta) = \frac{1}{\#\Theta}. \quad (2.12)$$

Proof. For any prior Π and any $\theta \in \Theta$ such that $\Pi(\theta) > 0$,

$$\begin{aligned} D(P_{\theta}^{\otimes n} \| P_{\Pi,n}) &= \int_{\mathcal{Y}^n} \log \left(\frac{dP_{\theta}^{\otimes n}(y)}{\sum_{\theta' \in \Theta} \Pi(\theta') dP_{\theta'}^{\otimes n}(y)} \right) dP_{\theta}^{\otimes n}(y) \\ &= - \int_{\mathcal{Y}^n} \log \left(\Pi(\theta) + \sum_{\theta' \neq \theta} \Pi(\theta') \frac{dP_{\theta'}^{\otimes n}(y)}{dP_{\theta}^{\otimes n}(y)} \right) dP_{\theta}^{\otimes n}(y) \\ &= - \log(\Pi(\theta)) - \int_{\mathcal{Y}^n} \log \left(1 + \sum_{\theta' \neq \theta} \frac{\Pi(\theta')}{\Pi(\theta)} \frac{dP_{\theta'}^{\otimes n}(y)}{dP_{\theta}^{\otimes n}(y)} \right) dP_{\theta}^{\otimes n}(y). \end{aligned} \quad (2.13)$$

Now, since Θ is finite and for all $\theta_1, \theta_2 \in \Theta$, $P_{\theta_1} \neq P_{\theta_2}$, there exists a weakly consistent estimator $\hat{\theta}_n$ for the model.

$$\begin{aligned}
& \int_{\mathcal{Y}^n} \log \left(1 + \sum_{\theta' \neq \theta} \frac{\Pi(\theta')}{\Pi(\theta)} \frac{dP_{\theta'}^{\otimes n}}{dP_{\theta}^{\otimes n}}(y) \right) dP_{\theta}^{\otimes n}(y) \\
&= \int_{\hat{\theta}_n = \theta} \log \left(1 + \sum_{\theta' \neq \theta} \frac{\Pi(\theta')}{\Pi(\theta)} \frac{dP_{\theta'}^{\otimes n}}{dP_{\theta}^{\otimes n}}(y) \right) dP_{\theta}^{\otimes n}(y) \\
&\quad + \frac{P_{\theta}^{\otimes n}(\hat{\theta}_n \neq \theta)}{P_{\theta}^{\otimes n}(\hat{\theta}_n \neq \theta)} \int_{\hat{\theta}_n \neq \theta} \log \left(1 + \sum_{\theta' \neq \theta} \frac{\Pi(\theta')}{\Pi(\theta)} \frac{dP_{\theta'}^{\otimes n}}{dP_{\theta}^{\otimes n}}(y) \right) dP_{\theta}^{\otimes n}(y) \\
&\leq \int_{\hat{\theta}_n = \theta} \left(\sum_{\theta' \neq \theta} \frac{\Pi(\theta')}{\Pi(\theta)} \frac{dP_{\theta'}^{\otimes n}}{dP_{\theta}^{\otimes n}}(y) \right) dP_{\theta}^{\otimes n}(y) \\
&\quad + P_{\theta}^{\otimes n}(\hat{\theta}_n \neq \theta) \log \left(1 + \frac{1}{P_{\theta}^{\otimes n}(\hat{\theta}_n \neq \theta)} \int_{\hat{\theta}_n \neq \theta} \sum_{\theta' \neq \theta} \frac{\Pi(\theta')}{\Pi(\theta)} \frac{dP_{\theta'}^{\otimes n}}{dP_{\theta}^{\otimes n}}(y) dP_{\theta}^{\otimes n}(y) \right) \\
&= \sum_{\theta' \neq \theta} \frac{\Pi(\theta')}{\Pi(\theta)} P_{\theta'}^{\otimes n}(\hat{\theta}_n = \theta) + P_{\theta}^{\otimes n}(\hat{\theta}_n \neq \theta) \log \left(1 + \sum_{\theta' \neq \theta} \frac{\Pi(\theta')}{\Pi(\theta)} \frac{P_{\theta'}^{\otimes n}(\hat{\theta}_n \neq \theta)}{P_{\theta}^{\otimes n}(\hat{\theta}_n \neq \theta)} \right) \\
&\leq \sum_{\theta' \neq \theta} \frac{\Pi(\theta')}{\Pi(\theta)} P_{\theta'}^{\otimes n}(\hat{\theta}_n = \theta) + P_{\theta}^{\otimes n}(\hat{\theta}_n \neq \theta) \log \left(1 + \sum_{\theta' \neq \theta} \frac{\Pi(\theta')}{\Pi(\theta)} \frac{1}{P_{\theta}^{\otimes n}(\hat{\theta}_n \neq \theta)} \right) \quad (2.14)
\end{aligned}$$

The first equality holds unless $P_{\theta}^{\otimes n}(\hat{\theta}_n \neq \theta) = 0$, but Equation (2.14) trivially holds in that case. The first inequality results from applying the Jensen inequality to the convex mapping $\log(1+x)$ and from the inequality $\log(1+x) \leq x$, which holds for all $x \geq 0$.

Since $\hat{\theta}_n$ is weakly consistent, we have both $\lim_{n \rightarrow +\infty} P_{\theta'}^{\otimes n}(\hat{\theta}_n = \theta) = 0$ for all $\theta' \neq \theta$ and $\lim_{n \rightarrow +\infty} P_{\theta}^{\otimes n}(\hat{\theta}_n \neq \theta) = 0$. Therefore

$$\lim_{n \rightarrow +\infty} D(P_{\theta}^{\otimes n} \| P_{\Pi, n}) = -\log(\Pi(\theta)) \quad (2.15)$$

and

$$\lim_{n \rightarrow +\infty} \sum_{\theta \in \Theta} D(P_{\theta}^{\otimes n} \| P_{\Pi, n}) \Pi(\theta) = -\sum_{\theta \in \Theta} \log(\Pi(\theta)) \Pi(\theta). \quad (2.16)$$

Because the logarithm is convex, for any prior Π ,

$$\sum_{\theta \in \Theta} \log \left(\frac{1}{\Pi(\theta)} \right) \Pi(\theta) \leq \log \left(\sum_{\theta \in \Theta} \frac{1}{\Pi(\theta)} \Pi(\theta) \right) = \log(\#\Theta). \quad (2.17)$$

This inequality becomes an equality if, and only if, $\Pi(\theta) = 1/\#\Theta$ for every $\theta \in \Theta$. □

Remark. The right member of (2.16) is the Shannon entropy of Π .

Interestingly, the reference prior here fits the Laplace [1814] insufficient reason principle, which consists in assuming all elementary events to be equiprobable. This is the first example of the reference prior naturally fitting earlier conceptions of noninformativity.

Although Definition 2.2 intuitively makes sense, it is too restrictive. The limit in (2.11) is often infinite, which means no reference prior exists. Nevertheless, it is sufficient for cases where \mathcal{Y} is a finite set.

2.3 Full definition of the reference prior

The definition of a reference prior below is a more restrictive version of the definition given by Bernardo [2005], because it requires such a prior to be proper.

Definition 2.4 (Reference prior: proper case). *A reference prior Π^* is an element of $\mathcal{P}(\Theta)$ such that for any other element Π of $\mathcal{P}(\Theta)$,*

$$\liminf_{n \rightarrow +\infty} \left[\int_{\Theta} D(P_{\theta}^{\otimes n} \| P_{\Pi^*, n}) d\Pi^*(\theta) - \int_{\Theta} D(P_{\theta}^{\otimes n} \| P_{\Pi, n}) d\Pi(\theta) \right] \geq 0. \quad (2.18)$$

Definition 2.4 is an extension of Definition 2.2. Any reference prior in the sense of Definition 2.2 is one in the sense of Definition 2.4.

Remark. If $(P_{\theta})_{\theta \in \Theta}$ is a dominated model, the following alternate definition is equivalent: a reference prior Π^* is an element of $\mathcal{P}(\Theta)$ such that for any other element Π of $\mathcal{P}(\Theta)$,

$$\liminf_{n \rightarrow +\infty} \left[\int_{\mathcal{Y}^n} D(\Pi^*(\cdot | y) \| \Pi^*) dP_{\Pi^*, n}(y) - \int_{\mathcal{Y}^n} D(\Pi(\cdot | y) \| \Pi) dP_{\Pi, n}(y) \right] \geq 0. \quad (2.19)$$

This definition is more in line with the first intuition presented in this chapter, which can now be refined. When $n \rightarrow +\infty$, the posterior distribution approaches perfect knowledge of the parameter. The Kullback-Leibler divergence between the posterior and prior can thus be seen as the difference between (almost) perfect knowledge and knowledge due to the prior. It is the knowledge “missing” from the prior [Berger et al., 2009].

One question that naturally arises is whether the reference prior, in case it exists, is unique. This part of the discussion is to our knowledge new, in the sense that uniqueness does not seem to have been studied in a general setting before. Bernardo [2005] and Berger et al. [2009] do not care about it, which is understandable, since two reference priors would necessarily be equally noninformative, so any of them could be chosen.

Lemma 2.5. *Let $t \in (0, 1)$ and let Π_t , Π_0 and Π_1 be proper priors such that $\Pi_t = (1-t)\Pi_0 + t\Pi_1$. For every positive integer n ,*

$$\int_{\Theta} d\Pi_1(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_t, n}) - \int_{\Theta} d\Pi_1(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_1, n}) = D(P_{\Pi_1, n} \| P_{\Pi_t, n}) \leq -\log(t). \quad (2.20)$$

Proof. The inequality is due to the fact that $P_{\Pi_t, n} = (1-t)P_{\Pi_0, n} + tP_{\Pi_1, n}$. So Π_1 -almost surely, $dP_{\Pi_1, n}/dP_{\Pi_t, n} \leq 1/t$, so

$$D(P_{\Pi_1, n} \| P_{\Pi_t, n}) = \int_{\mathcal{Y}^n} dP_{\Pi_1, n}(y) \log \left(\frac{dP_{\Pi_1, n}(y)}{dP_{\Pi_t, n}(y)} \right) \leq \log \left(\frac{1}{t} \right). \quad (2.21)$$

To obtain the equality, notice that

$$\begin{aligned} & \int_{\Theta} d\Pi_1(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_t, n}) \\ &= \int_{\Theta} d\Pi_1(\theta) \int_{\mathcal{Y}^n} dP_{\theta}^{\otimes n}(\mathbf{y}) \log \left(\frac{dP_{\theta}^{\otimes n}(\mathbf{y})}{dP_{\Pi_t, n}(\mathbf{y})} \right) \\ &= \int_{\Theta} d\Pi_1(\theta) \int_{\mathcal{Y}^n} dP_{\theta}^{\otimes n}(\mathbf{y}) \log \left(\frac{dP_{\theta}^{\otimes n}(\mathbf{y})}{dP_{\Pi_1, n}(\mathbf{y})} \right) + \int_{\Theta} d\Pi_1(\theta) \int_{\mathcal{Y}^n} dP_{\theta}^{\otimes n}(\mathbf{y}) \log \left(\frac{dP_{\Pi_1, n}(\mathbf{y})}{dP_{\Pi_t, n}(\mathbf{y})} \right) \\ &= \int_{\Theta} d\Pi_1(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_1, n}) + D(P_{\Pi_1, n} \| P_{\Pi_t, n}). \end{aligned} \quad (2.22)$$

The last step holds due to the Fubini-Lebesgue theorem, which can be used because

$$\begin{aligned} \int_{\Theta} d\Pi_1(\theta) \int_{\mathcal{Y}^n} dP_{\theta}^{\otimes n}(\mathbf{y}) \left| \log \left(\frac{dP_{\Pi_1, n}(\mathbf{y})}{dP_{\Pi_t, n}(\mathbf{y})} \right) \right| &\leq \int_{\Theta} d\Pi_1(\theta) \int_{\mathcal{Y}^n} dP_{\theta}^{\otimes n}(\mathbf{y}) \log \left(\frac{dP_{\Pi_1, n}(\mathbf{y})}{dP_{\Pi_t, n}(\mathbf{y})} \right) + 1 \\ &\leq \int_{\Theta} d\Pi_1(\theta) \int_{\mathcal{Y}^n} dP_{\theta}^{\otimes n}(\mathbf{y}) \log \left(\frac{1}{t} \right) + 1 \\ &= -\log(t) + 1 < +\infty. \end{aligned} \quad (2.23)$$

□

Lemma 2.6. *For any $t \in (0, 1)$, for any proper priors Π_t , Π_0 and Π_1 such that $\Pi_t = (1-t)\Pi_0 + t\Pi_1$,*

$$\begin{aligned} \int_{\Theta} d\Pi_t(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_t, n}) &= (1-t) \int_{\Theta} d\Pi_0(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_0, n}) + t \int_{\Theta} d\Pi_1(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_1, n}) \\ &\quad + (1-t) D(P_{\Pi_0, n} \| P_{\Pi_t, n}) + t D(P_{\Pi_1, n} \| P_{\Pi_t, n}). \end{aligned} \quad (2.24)$$

Proof. The result follows from this decomposition:

$$\int_{\Theta} d\Pi(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_t, n}) = (1-t) \int_{\Theta} d\Pi_0(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_t, n}) + t \int_{\Theta} d\Pi_1(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_t, n}). \quad (2.25)$$

Applying Lemma 2.5 to both Π_0 and Π_1 yields the result. □

Lemma 2.7. *For any $t \in (0, 1)$, for any proper priors Π_t , Π_0 and Π_1 such that $\Pi_t = (1-t)\Pi_0 + t\Pi_1$, the sequences $(D(P_{\Pi_0, n} \| P_{\Pi_t, n}))_{n \in \mathbb{Z}_+}$ and $(D(P_{\Pi_1, n} \| P_{\Pi_t, n}))_{n \in \mathbb{Z}_+}$ are nondecreasing.*

Proof. We only prove that $(D(P_{\Pi_1, n} \| P_{\Pi_t, n}))_{n \in \mathbb{Z}_+}$ is nondecreasing because the proof that $(D(P_{\Pi_0, n} \| P_{\Pi_t, n}))_{n \in \mathbb{Z}_+}$ is nondecreasing is similar.

Recall that for every positive integer n , $P_{\Pi_t, n} = (1-t)P_{\Pi_0, n} + tP_{\Pi_1, n}$ and therefore $P_{\Pi_0, n}$ -almost surely, $dP_{\Pi_1, n}/dP_{\Pi_t, n} \leq 1/t$.

$$\begin{aligned} D(P_{\Pi_1, n} \| P_{\Pi_t, n}) &= \int_{\mathcal{Y}^n} dP_{\Pi_1, n}(\mathbf{y}) \log \left(\frac{dP_{\Pi_1, n}(\mathbf{y})}{dP_{\Pi_t, n}(\mathbf{y})} \right) \\ &\leq \int_{\mathcal{Y}^n} dP_{\Pi_1, n}(\mathbf{y}) \log \left(\frac{1}{1-t} \right) = -\log(t) < +\infty. \end{aligned} \quad (2.26)$$

Therefore

$$\int_{\mathcal{Y}^n} dP_{\Pi_1, n}(\mathbf{y}) \left| \log \left(\frac{dP_{\Pi_1, n}(\mathbf{y})}{dP_{\Pi_t, n}(\mathbf{y})} \right) \right| \leq -\log(t) + 1 < +\infty. \quad (2.27)$$

Now, define $\Pi_{0, n}(\cdot | \mathbf{y})$ (resp. $\Pi_{1, n}(\cdot | \mathbf{y})$, $\Pi_{t, n}(\cdot | \mathbf{y})$) as the posterior distribution resulting from the prior Π_1 (resp. Π_0 , Π_t) after having made n independent observations (\mathbf{y}) . Then define the ‘‘conditional’’ predictive distribution $P_{\Pi_1|n}^{\mathbf{y}} = \int_{\Theta} d\Pi_1, n(\theta | \mathbf{y}) P_{\theta}$ (resp. $P_{\Pi_0|n}^{\mathbf{y}} = \int_{\Theta} d\Pi_0, n(\theta | \mathbf{y}) P_{\theta}$, $P_{\Pi_t|n}^{\mathbf{y}} = \int_{\Theta} d\Pi_t, n(\theta | \mathbf{y}) P_{\theta}$). Applying the Fubini-Lebesgue theorem,

$$\begin{aligned}
D(P_{\Pi_1, n+1} \| P_{\Pi_t, n+1}) &= \int_{\mathcal{Y}^n} dP_{\Pi_1, n}(\mathbf{y}) \int_{\mathcal{Y}} dP_{\Pi_1|n}^{\mathbf{y}}(y') \log \left(\frac{dP_{\Pi_1, n+1}}{dP_{\Pi_t, n+1}}(\mathbf{y}, y') \right) \\
&= \int_{\mathcal{Y}^n} dP_{\Pi_1, n}(\mathbf{y}) \int_{\mathcal{Y}} dP_{\Pi_1|n}^{\mathbf{y}}(y') \log \left(\frac{dP_{\Pi_1, n}}{dP_{\Pi_t, n}}(\mathbf{y}) \frac{dP_{\Pi_1|n}^{\mathbf{y}}}{dP_{\Pi_t|n}^{\mathbf{y}}}(y') \right) \\
&= \int_{\mathcal{Y}^n} dP_{\Pi_1, n}(\mathbf{y}) \log \left(\frac{dP_{\Pi_1, n}}{dP_{\Pi_t, n}}(\mathbf{y}) \right) \\
&\quad + \int_{\mathcal{Y}^n} dP_{\Pi_1, n}(\mathbf{y}) \int_{\mathcal{Y}} dP_{\Pi_1|n}^{\mathbf{y}}(y') \log \left(\frac{dP_{\Pi_1|n}^{\mathbf{y}}}{dP_{\Pi_t|n}^{\mathbf{y}}}(y') \right) \\
&= D(P_{\Pi_1, n} \| P_{\Pi_t, n}) + \int_{\mathcal{Y}^n} dP_{\Pi_1, n}(\mathbf{y}) D(P_{\Pi_1|n}^{\mathbf{y}} \| P_{\Pi_t|n}^{\mathbf{y}}) \\
&\geq D(P_{\Pi_1, n} \| P_{\Pi_t, n}). \tag{2.28}
\end{aligned}$$

The second equality holds because $P_{\Pi_t, n+1} = (1-t)P_{\Pi_1, n+1} + tP_{\Pi_t, n+1}$ and therefore

- $P_{\Pi_1, n}$ -almost surely, $dP_{\Pi_1, n}/dP_{\Pi_t, n} \leq 1/t$;
- for Π_1, n -almost any $\mathbf{y} \in \mathcal{Y}^n$, $P_{\Pi_1|n}^{\mathbf{y}}$ -almost surely, $dP_{\Pi_1|n}^{\mathbf{y}}/dP_{\Pi_t|n}^{\mathbf{y}} \leq 1/t$.

□

Lemma 2.8. *For any proper prior Π , the mapping $i_{\Pi} : \mathcal{P}(\Theta) \rightarrow [0, +\infty]$ defined by*

$$i_{\Pi}(\Pi') := \liminf_{n \rightarrow +\infty} \left[\int_{\Theta} D(P_{\theta}^{\otimes n} \| P_{\Pi', n}) d\Pi'(\theta) - \int_{\Theta} D(P_{\theta}^{\otimes n} \| P_{\Pi, n}) d\Pi(\theta) \right] \tag{2.29}$$

is essentially strictly convex in the following sense. For all $\Pi_0, \Pi_1 \in \mathcal{P}(\Theta)$ such that $i_{\Pi}(\Pi_0)$ and $i_{\Pi}(\Pi_1)$ are finite,

- *either for all $t \in (0, 1)$, $i_{\Pi}((1-t)\Pi_0 + t\Pi_1) > (1-t)i_{\Pi}(\Pi_0) + ti_{\Pi}(\Pi_1)$;*
- *or for all $n \in \mathbb{Z}_+$, $P_{\Pi_0, n} = P_{\Pi_1, n}$.*

Proof. Define $\Pi_t := (1-t)\Pi_0 + t\Pi_1$. Lemmas 2.6 and 2.7 imply that for all $m \in \mathbb{Z}_+$,

$$\begin{aligned}
&\liminf_{n \rightarrow \infty} \left[\int_{\Theta} d\Pi_t(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_t, n}) \right. \\
&\quad \left. - (1-t) \int_{\Theta} d\Pi_0(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_0, n}) - t \int_{\Theta} d\Pi_1(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_1, n}) \right] \\
&\geq (1-t)D(P_{\Pi_0, m} \| P_{\Pi_t, m}) + tD(P_{\Pi_1, m} \| P_{\Pi_t, m}). \tag{2.30}
\end{aligned}$$

For any proper prior Π such that $i_{\Pi}(\Pi_0)$ and $i_{\Pi}(\Pi_1)$ are finite,

$$\begin{aligned}
&\liminf_{n \rightarrow \infty} \int_{\Theta} d\Pi_t(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_t, n}) \\
&\quad - (1-t) \int_{\Theta} d\Pi_0(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_0, n}) - t \int_{\Theta} d\Pi_1(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_1, n}) \\
&= \liminf_{n \rightarrow \infty} \left[\int_{\Theta} d\Pi_t(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_t, n}) - \int_{\Theta} d\Pi(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi, n}) \right. \\
&\quad \left. - (1-t) \left(\int_{\Theta} d\Pi_0(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_0, n}) - \int_{\Theta} d\Pi(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi, n}) \right) \right. \\
&\quad \left. - t \left(\int_{\Theta} d\Pi_1(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_1, n}) - \int_{\Theta} d\Pi(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi, n}) \right) \right] \tag{2.31}
\end{aligned}$$

$$\begin{aligned}
& \liminf_{n \rightarrow \infty} \int_{\Theta} d\Pi_t(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_t, n}) \\
& \quad - (1-t) \int_{\Theta} d\Pi_0(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_0, n}) - t \int_{\Theta} d\Pi_1(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_1, n}) \\
\leq & \liminf_{n \rightarrow \infty} \left[\int_{\Theta} d\Pi_t(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_t, n}) - \int_{\Theta} d\Pi(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi, n}) \right] \\
& + \limsup_{n \rightarrow +\infty} \left[-(1-t) \left(\int_{\Theta} d\Pi_0(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_0, n}) - \int_{\Theta} d\Pi(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi, n}) \right) \right] \\
& + \limsup_{n \rightarrow +\infty} \left[-t \left(\int_{\Theta} d\Pi_1(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_1, n}) - \int_{\Theta} d\Pi(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi, n}) \right) \right] \\
= & \liminf_{n \rightarrow \infty} \left[\int_{\Theta} d\Pi_t(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_t, n}) - \int_{\Theta} d\Pi(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi, n}) \right] \\
& - (1-t) \liminf_{n \rightarrow \infty} \left[\int_{\Theta} d\Pi_0(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_0, n}) - \int_{\Theta} d\Pi(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi, n}) \right] \\
& - t \liminf_{n \rightarrow \infty} \left[\int_{\Theta} d\Pi_1(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi_1, n}) - \int_{\Theta} d\Pi(\theta) D(P_{\theta}^{\otimes n} \| P_{\Pi, n}) \right] \\
= & i_{\Pi}(\Pi_t) - (1-t)i_{\Pi}(\Pi_0) - ti_{\Pi}(\Pi_1). \tag{2.32}
\end{aligned}$$

Combining Equations (2.30) and (2.32), we obtain that for all proper priors Π , Π_0 and Π_1 such that both $i_{\Pi}(\Pi_0)$ and $i_{\Pi}(\Pi_1)$ are finite and for every positive integer n ,

$$i_{\Pi}(\Pi_t) - (1-t)i_{\Pi}(\Pi_0) - ti_{\Pi}(\Pi_1) \geq (1-t)D(P_{\Pi_0, n} \| P_{\Pi_t, n}) + tD(P_{\Pi_1, n} \| P_{\Pi_t, n}). \tag{2.33}$$

Unless $P_{\Pi_0, n} = P_{\Pi_1, n}$ for every positive integer n , there exists a positive integer m such that

$$(1-t)D(P_{\Pi_0, m} \| P_{\Pi_t, m}) + tD(P_{\Pi_1, m} \| P_{\Pi_t, m}) > 0. \tag{2.34}$$

Therefore

$$i_{\Pi}(\Pi_t) - (1-t)i_{\Pi}(\Pi_0) - ti_{\Pi}(\Pi_1) > 0. \tag{2.35}$$

□

Let us make the following assumption:

Assumption 1. For every pair of distinct proper priors Π_0 and Π_1 on Θ , there exists a positive integer n such that $P_{\Pi_0, n} \neq P_{\Pi_1, n}$.

Proposition 2.9. Under Assumption 1, if there exists a reference prior in the sense of Definition 2.4, then it is unique.

Proof. Let Π_0^* and Π_1^* be reference priors. Then $i_{\Pi_0^*}(\Pi_1^*) = 0$. Besides, we also trivially have $i_{\Pi_0^*}(\Pi_0^*) = 0$. By Lemma 2.8,

- either for all $t \in (0, 1)$, $i_{\Pi_0^*}((1-t)\Pi_0^* + t\Pi_1^*) > 0$, which is impossible since it would contradict Definition 2.4;
- or for every positive integer n , $P_{\Pi_0^*, n} = P_{\Pi_1^*, n}$.

By Assumption 1, this implies that $\Pi_0^* = \Pi_1^*$. □

If the model $(P_\theta)_{\theta \in \Theta}$ is identifiable ($P_\theta = P_{\theta'} \Rightarrow \theta = \theta'$) and under some topological conditions, a de Finetti-type theorem may apply and make Assumption 1 hold (see Ressel [1985]). Sometimes more elementary arguments apply:

Proposition 2.10. *If there exists a weakly consistent estimator of θ , then Assumption 1 holds.*

Proof. For any subset $\Theta' \subset \Theta$ and any $\theta \in \Theta$, define $\text{dist}(\theta, \Theta') := \inf\{\text{dist}(\theta, \theta') : \theta' \in \Theta'\}$.

Let the sequence of mappings $\hat{\theta}_n : \mathcal{Y}^n \rightarrow \Theta$ ($n \in \mathbb{N}$) be a weakly consistent estimator. This means that for any $\theta \in \Theta$, for any $\epsilon > 0$, $\lim_{n \rightarrow +\infty} P_\theta^{\otimes n}(\text{dist}(\hat{\theta}_n, \theta) < \epsilon) = 1$. It implies that for any open set U , for any $\theta \in U$, $\lim_{n \rightarrow +\infty} P_\theta^{\otimes n}(\hat{\theta}_n \in U) = 1$ and that for any $\theta \in \Theta$ such that $\text{dist}(\theta, U) > 0$, $\lim_{n \rightarrow +\infty} P_\theta^{\otimes n}(\hat{\theta}_n \in U) = 0$.

Let U be an open subset of Θ . For all $\epsilon > 0$, define $U_\epsilon := \{\theta \in U : B(\theta, \epsilon) \subset U\}$, where $B(\theta, \epsilon)$ is the open ball centered around θ and with radius ϵ : $B(\theta, \epsilon) := \{\theta' \in \Theta : \text{dist}(\theta, \theta') < \epsilon\}$. U_ϵ is an open, possibly empty, set.

Now, let Π_0 and Π_1 be two proper prior distributions on Θ . Applying the dominated convergence theorem,

$$\Pi_0(U_\epsilon) = \int_{U_\epsilon} d\Pi_0(\theta) \lim_{n \rightarrow +\infty} P_\theta^{\otimes n}(\hat{\theta}_n \in U_\epsilon) = \lim_{n \rightarrow +\infty} \int_{U_\epsilon} d\Pi_0(\theta) P_\theta^{\otimes n}(\hat{\theta}_n \in U_\epsilon). \quad (2.36)$$

For any $n \in \mathbb{N}$,

$$\int_{\Theta} d\Pi_0(\theta) P_\theta^{\otimes n}(\hat{\theta}_n \in U_\epsilon) \leq \Pi_0(U) + \int_{\Theta \setminus U} d\Pi_0(\theta) P_\theta^{\otimes n}(\hat{\theta}_n \in U_\epsilon). \quad (2.37)$$

Again, applying the dominated convergence theorem,

$$\lim_{n \rightarrow +\infty} \int_{\Theta \setminus U} d\Pi_0(\theta) P_\theta^{\otimes n}(\hat{\theta}_n \in U_\epsilon) = \int_{\Theta \setminus U} d\Pi_0(\theta) \lim_{n \rightarrow +\infty} P_\theta^{\otimes n}(\hat{\theta}_n \in U_\epsilon) = 0. \quad (2.38)$$

Gathering the two equations above,

$$\limsup_{n \rightarrow +\infty} \int_{\Theta} d\Pi_0(\theta) P_\theta^{\otimes n}(\hat{\theta}_n \in U_\epsilon) \leq \Pi_0(U). \quad (2.39)$$

Combining this with Equation (2.36),

$$\Pi_0(U_\epsilon) \leq \limsup_{n \rightarrow +\infty} P_{\Pi_0, n}(\hat{\theta}_n \in U_\epsilon) \leq \Pi_0(U). \quad (2.40)$$

Let $(\epsilon_m)_{m \in \mathbb{Z}_+}$ be a decreasing sequence of positive real numbers with null limit. For example, $\epsilon_m = 1/m$. Then $U = \bigcup_{m \in \mathbb{Z}_+} U_{\epsilon_m}$. Since $(U_{\epsilon_m})_{m \in \mathbb{Z}_+}$ is an increasing sequence of open (hence measurable) sets, $\Pi_0(U) = \lim_{m \rightarrow +\infty} \Pi_0(U_{\epsilon_m})$.

In the end,

$$\Pi_0(U) = \lim_{m \rightarrow +\infty} \limsup_{n \rightarrow \infty} P_{\Pi_0, n}(\hat{\theta}_n \in U_{\epsilon_m}). \quad (2.41)$$

Similarly, we have

$$\Pi_1(U) = \lim_{m \rightarrow +\infty} \limsup_{n \rightarrow \infty} P_{\Pi_1, n}(\hat{\theta}_n \in U_{\epsilon_m}). \quad (2.42)$$

So, if for every positive integer n $P_{\Pi_0, n} = P_{\Pi_1, n}$, then for every open set U , $\Pi_0(U) = \Pi_1(U)$. Given a topology is a π -system that generates its Borel σ -algebra, this implies that $\Pi_0 = \Pi_1$.

As an aside, notice that a similar argument yields

$$\Pi_0(U) = \lim_{m \rightarrow +\infty} \liminf_{n \rightarrow \infty} P_{\Pi_0, n}(\hat{\theta}_n \in U_{\epsilon_m}) \quad (2.43)$$

and

$$\Pi_1(U) = \lim_{m \rightarrow +\infty} \liminf_{n \rightarrow \infty} P_{\Pi_1, n}(\hat{\theta}_n \in U_{\epsilon_m}). \quad (2.44)$$

□

The above served to prove that under a reasonable assumption – Assumption 1 – a proper reference prior is unique. Definition 2.13 below extends the definition of a reference prior in order to allow for improper reference priors. To ensure a similar uniqueness result holds with this extended definition, some preparatory work is needed.

First, a new assumption is needed.

Assumption 2. *If Π_0 and Π_1 are proper priors on Θ such that there exists an open subset $U \subset \Theta$ with $\Pi_0(U) = 0$ and $\Pi_1(U) = 1$, then there exists a sequence of sets $(T_n)_{n \in \mathbb{Z}_+}$ with the following properties:*

1. *For every positive integer n , $T_n \in \mathbb{Y}^{\otimes n}$.*
2. *$\lim_{n \rightarrow \infty} P_{\Pi_0, n}(T_n) = 0$ and $\lim_{n \rightarrow \infty} P_{\Pi_1, n}(T_n) = 1$.*

We later show that the existence of a weakly consistent estimator for the model $(P_\theta)_{\theta \in \Theta}$ implies Assumption 2.

Proposition 2.11. *Under Assumption 2, if a proper prior Π^* is a reference prior in the sense of Definition 2.4, then for any open subset U of Θ such that $\Pi^*(U) > 0$, the renormalized restriction of Π^* to U is a reference prior on U in the sense of Definition 2.4.*

Proof. Let Π^* be a proper reference prior in the sense of Definition 2.4. Let U be an open subset of Θ such that $\Pi^*(U) > 0$. Define $t := \Pi^*(U)$.

If $t = 1$, then there is nothing to prove, so assume $t < 1$.

Define Π_0^* as the normalized restriction of Π^* to $\Theta \setminus U$ and Π_1^* as the normalized restriction of Π^* to U . Therefore $\Pi^* = (1 - t)\Pi_0^* + t\Pi_1^*$.

Now let Π_0 be any proper prior with support on $\Theta \setminus U$ and Π_1 be any proper prior with support on U . Define $\Pi_t := (1 - t)\Pi_0 + t\Pi_1$.

Using Assumption 2, there exists a sequence $(T_n)_{n \in \mathbb{Z}_+}$ that has the properties 1. and 2.

For every positive integer n ,

$$\begin{aligned}
D(P_{\Pi_1, n} || P_{\Pi_t, n}) &= - \int_{\mathcal{Y}^n} dP_{\Pi_1, n}(\mathbf{y}) \log \left(\frac{dP_{\Pi_t, n}(\mathbf{y})}{dP_{\Pi_1, n}(\mathbf{y})} \right) \\
&= - \int_{\mathcal{Y}^n} dP_{\Pi_1, n}(\mathbf{y}) \log \left((1-t) \frac{dP_{\Pi_0, n}(\mathbf{y})}{dP_{\Pi_1, n}(\mathbf{y})} + t \right) \\
&= - \log(t) - \int_{\mathcal{Y}^n} dP_{\Pi_1, n}(\mathbf{y}) \log \left(1 + \frac{1-t}{t} \frac{dP_{\Pi_0, n}(\mathbf{y})}{dP_{\Pi_1, n}(\mathbf{y})} \right). \tag{2.45}
\end{aligned}$$

$$\begin{aligned}
&\int_{\mathcal{Y}^n} dP_{\Pi_1, n}(\mathbf{y}) \log \left(1 + \frac{1-t}{t} \frac{dP_{\Pi_0, n}(\mathbf{y})}{dP_{\Pi_1, n}(\mathbf{y})} \right) \\
&= \int_{T_n} dP_{\Pi_1, n}(\mathbf{y}) \log \left(1 + \frac{1-t}{t} \frac{dP_{\Pi_0, n}(\mathbf{y})}{dP_{\Pi_1, n}(\mathbf{y})} \right) + \int_{\mathcal{Y}^n \setminus T_n} dP_{\Pi_1, n}(\mathbf{y}) \log \left(1 + \frac{1-t}{t} \frac{dP_{\Pi_0, n}(\mathbf{y})}{dP_{\Pi_1, n}(\mathbf{y})} \right) \\
&\leq \int_{T_n} dP_{\Pi_1, n}(\mathbf{y}) \frac{1-t}{t} \frac{dP_{\Pi_0, n}(\mathbf{y})}{dP_{\Pi_1, n}(\mathbf{y})} + \int_{\mathcal{Y}^n \setminus T_n} dP_{\Pi_1, n}(\mathbf{y}) \log \left(1 + \frac{1-t}{t} \frac{dP_{\Pi_0, n}(\mathbf{y})}{dP_{\Pi_1, n}(\mathbf{y})} \right) \\
&= \frac{1-t}{t} P_{\Pi_0, n}(T_n) + \frac{P_{\Pi_1, n}(\mathcal{Y}^n \setminus T_n)}{P_{\Pi_1, n}(\mathcal{Y}^n \setminus T_n)} \int_{\mathcal{Y}^n \setminus T_n} dP_{\Pi_1, n}(\mathbf{y}) \log \left(1 + \frac{1-t}{t} \frac{dP_{\Pi_0, n}(\mathbf{y})}{dP_{\Pi_1, n}(\mathbf{y})} \right) \\
&\leq \frac{1-t}{t} P_{\Pi_0, n}(T_n) + P_{\Pi_1, n}(\mathcal{Y}^n \setminus T_n) \log \left(1 + \frac{1-t}{t} \frac{P_{\Pi_0, n}(\mathcal{Y}^n \setminus T_n)}{P_{\Pi_1, n}(\mathcal{Y}^n \setminus T_n)} \right) \\
&\leq \frac{1-t}{t} P_{\Pi_0, n}(T_n) + P_{\Pi_1, n}(\mathcal{Y}^n \setminus T_n) \log \left(1 + \frac{1-t}{t} \frac{1}{P_{\Pi_1, n}(\mathcal{Y}^n \setminus T_n)} \right). \tag{2.46}
\end{aligned}$$

In the computation above, the first inequality holds because for any nonnegative real number x , $\log(1+x) \leq x$. The second equality only makes sense if $P_{\Pi_1, n}(\mathcal{Y}^n \setminus T_n) \neq 0$, but in case $P_{\Pi_1, n}(\mathcal{Y}^n \setminus T_n) = 0$, Equation 2.46 holds a fortiori. The second inequality results from Jensen's inequality applied to the concave logarithm function.

Due to Assumption 2, $\lim_{n \rightarrow +\infty} P_{\Pi_0, n}(T_n) = \lim_{n \rightarrow +\infty} P_{\Pi_1, n}(\mathcal{Y}^n \setminus T_n) = 0$. Therefore, combining Equations (2.45) and (2.46) yields

$$\lim_{n \rightarrow +\infty} D(P_{\Pi_1, n} || P_{\Pi_t, n}) = -\log(t). \tag{2.47}$$

A similar reasoning yields

$$\lim_{n \rightarrow +\infty} D(P_{\Pi_0, n} || P_{\Pi_t, n}) = -\log(1-t). \tag{2.48}$$

These two results do not actually depend on Π_0 and Π_1 but only on the fact that $\Pi_0(U) = 0$ and $\Pi_1(U) = 1$. Therefore they also hold for Π_0^* and Π_1^* .

Let us take $\Pi_0 = \Pi_0^*$.

Lemma 2.6, applied to Π^* and Π_t , implies that

$$i_{\Pi_t}(\Pi^*) = t i_{\Pi_1}(\Pi_1^*). \tag{2.49}$$

Because Π^* is a reference prior, this implies that $i_{\Pi_1}(\Pi_1^*) \geq 0$. Given this holds for any proper prior Π_1 such that $\Pi_1(U) = 1$, Π_1^* is a reference prior on U . \square

Intuitively, Assumption 2 states that for radically different proper priors (that do not even have the same support), prediction should tend to become radically different too when the number of observations increases. We only require this for “reasonable” supports, that is

supports that can be expressed in terms of open and closed sets. This is consistent with the general idea that the topology of Θ is important: P_θ and $P_{\theta'}$ should be close if $\text{dist}(\theta, \theta')$ is small.

Proposition 2.12. *If there exists a weakly consistent estimator of θ , then Assumption 2 holds.*

Proof. Let the sequence of mappings $\hat{\theta}_n : \mathcal{Y}^n \rightarrow \Theta$ ($n \in \mathbb{N}$) be a weakly consistent estimator. This means that for any $\theta \in \Theta$, for any $\epsilon > 0$, $\lim_{n \rightarrow +\infty} P_\theta^{\otimes n}(\text{dist}(\hat{\theta}_n, \theta) < \epsilon) = 1$.

Let Π_0 and Π_1 be two proper priors on Θ such that there exists an open subset $U \subset \Theta$ with $\Pi_0(U) = 0$ and $\Pi_1(U) = 1$. Notations are the same as in the proof of Proposition 2.10.

Define $N(0) := 0$.

For every integer $m = 1, 2, \dots$, do the following:

1. Choose an integer $N_0(m) > N(m-1)$ such that for every integer $n \geq N_0(m)$, $P_{\Pi_0, n}(\hat{\theta}_n \in U_{\epsilon_m}) < \limsup_{n' \rightarrow \infty} P_{\Pi_0, n'}(\hat{\theta}_{n'} \in U_{1/m}) + 1/m$.
2. Choose an integer $N_1(m) > N(m-1)$ such that for every integer $n \geq N_1(m)$, $P_{\Pi_1, n}(\hat{\theta}_n \in U_{\epsilon_m}) > \liminf_{n' \rightarrow \infty} P_{\Pi_1, n'}(\hat{\theta}_{n'} \in U_{1/m}) - 1/m$.
3. Define $N(m) := \max(N_0(m), N_1(m))$ and $T_{N(m)}$ as the event $\{\hat{\theta}_{N(m)} \in U_{1/m}\} \in \mathbb{Y}^{\otimes N(m)}$.

For every positive integer n , there exists a nonnegative integer m such that $N(m) \leq n < N(m+1)$. Define $T_n := T_{N(m)} \times \mathcal{Y}^{n-N(m)}$. We have:

$$P_{\Pi_0, n}(T_n) = P_{\Pi_0, N(m)}(T_{N(m)}) < \limsup_{n' \rightarrow \infty} P_{\Pi_0, n'}(\hat{\theta}_{n'} \in U_{1/m}) + 1/m; \quad (2.50)$$

$$P_{\Pi_1, n}(T_n) = P_{\Pi_1, N(m)}(T_{N(m)}) > \liminf_{n' \rightarrow \infty} P_{\Pi_1, n'}(\hat{\theta}_{n'} \in U_{1/m}) - 1/m. \quad (2.51)$$

Following Equations (2.41) and (2.44),

$$0 = \Pi_0(U) = \lim_{m \rightarrow +\infty} \limsup_{n \rightarrow \infty} P_{\Pi_0, n}(\hat{\theta}_n \in U_{1/m}); \quad (2.52)$$

$$1 = \Pi_1(U) = \lim_{m \rightarrow +\infty} \liminf_{n \rightarrow \infty} P_{\Pi_1, n}(\hat{\theta}_n \in U_{1/m}). \quad (2.53)$$

Therefore $\lim_{n \rightarrow +\infty} P_{\Pi_0, n}(T_n) = 0$ and $\lim_{n \rightarrow +\infty} P_{\Pi_1, n}(T_n) = 1$. □

The requirement that a reference prior should be proper (that is belong to $\mathcal{P}(\Theta)$) is counter-intuitive. On the contrary, one would rather expect it to be improper in some cases since it is supposed to have the least influence possible on inference [Berger and Bernardo, 1992]. The propriety requirement can be lifted in the following way:

Definition 2.13 (Reference prior: final definition). *A reference prior Π^* is a measure such that there exists an increasing sequence of open subsets $(U_n)_{n \in \mathbb{N}}$ ($U_0 \subset U_1 \subset U_2 \subset \dots$) such that $\bigcup_{n \in \mathbb{N}} U_n = \Theta$ and for every nonnegative integer n , $\Pi^*(U_n) < +\infty$ and which has the following property. For every nonnegative integer n , the probability distribution Π_n^* with support on U_n defined by $\Pi_n^*(A) = \Pi^*(A)/\Pi^*(U_n)$ for every measurable set $A \subset U_n$ is a reference prior on U_n in the sense of Definition 2.4).*

Remark. This definition is slightly different from Definition 6 in Bernardo [2005] in that Bernardo only requires the increasing sequence to consist of measurable sets and not necessarily open sets. This restriction, which is of little practical relevance, is used to obtain uniqueness results.

Proposition 2.14. *Under Assumptions 1 and 2, if a reference prior in the sense of Definition 2.13 exists, it is unique up to a multiplicative constant.*

Proof. Let Π_0^* and Π_1^* be two reference priors in the sense of Definition 2.13. So there exists an increasing sequence of open subsets $(U_m^0)_{m \in \mathbb{N}}$ such that for each nonnegative integer m , the renormalized restriction of Π_0^* to U_m^0 is a reference prior on U_m^0 . There also exists an increasing sequence of open subsets $(U_m^1)_{m \in \mathbb{N}}$ which plays the same role with respect to Π_1^* .

Let us define for every nonnegative integer m the open subset of Θ $V_m := U_m^0 \cap U_m^1$. Then $(V_m)_{m \in \mathbb{N}}$ is an increasing sequence of open subsets of Θ . Moreover, Assumption 2 and Proposition 2.11 imply that for every nonnegative integer m , both the renormalized restriction of Π_0^* to V_m and the renormalized restriction of Π_1^* to V_m are reference priors on V_m in the sense of Definition 2.4. So Assumption 1 and Proposition 2.9 imply that both renormalized restrictions to V_m are equal. Given this holds for every V_m ($m \in \mathbb{N}$), Π_0^* and Π_1^* are proportional to one another. □

This ends the original discussion about uniqueness.

2.4 Regular continuous case

In this section, we consider a particular case of great theoretical and practical importance. Everything here apart from Lemma 2.17 and Theorem 2.18 is taken from Clarke and Barron [1994].

Consider the following conditions.

Condition 0. Θ is a nonempty open subset of \mathbb{R}^d (relative to the Euclidean norm $\|\cdot\|$) whose boundary has null d -dimensional Lebesgue measure. For any $\theta \in \Theta$, P_θ be absolutely continuous with respect to a measure μ . Let p_θ be the Radon-Nykodym derivative; for μ -almost any y , the mapping $\theta \mapsto p_\theta(y)$ is twice continuously differentiable over Θ .

Condition 1. There exists $\epsilon > 0$ such that for all $\theta \in \Theta$, there exists $\delta(\theta) > 0$ such that for any integers $j, k \in \llbracket 1, d \rrbracket$, both

$$f(\theta) := \int_{\mathcal{Y}} \sup_{\theta' \in \Theta: \|\theta' - \theta\| < \delta(\theta)} \left| \frac{\partial^2}{\partial \theta_j' \partial \theta_k'} \log p_{\theta'}(y) \right|^2 p_{\theta'}(y) d\mu(y)$$

$$\text{and } g_\epsilon(\theta) := \int_{\mathcal{Y}} \left| \frac{\partial}{\partial \theta_j} \log p_\theta(y) \right|^{2+\epsilon} p_\theta(y) d\mu(y)$$

are finite. Moreover, f and g_ϵ are continuous functions of $\theta \in \Theta$.

Condition 2. For all $\theta \in \Theta$, for any integers $j, k \in \llbracket 1, d \rrbracket$,

$$\int_{\mathcal{Y}} \frac{\partial^2}{\partial \theta_j \partial \theta_k} p_\theta(y) d\mu(y) = 0.$$

Condition 3. For all $\theta, \theta' \in \Theta$, $P_\theta \neq P_{\theta'}$.

Under these conditions, the Jeffreys-rule prior is well defined on Θ .

Definition 2.15. *The Jeffreys-rule prior is the prior with density with respect to the Lebesgue measure proportional to the square root of the determinant of the Fisher information matrix: $|\mathcal{I}(\theta)|^{1/2}$.*

Remark. The Jeffreys-rule prior is traditionally viewed as noninformative because of its invariance by reparametrization and its reliance on Fisher information [Robert, 2007]. The Fisher information matrix, which is the variance of the score function and is linked to the curvature of the likelihood can be seen as representing the local discriminating power of the data [Robert et al., 2009].

The following theorem by Clarke and Barron [1994] concerns the asymptotic behavior of the noninformativity criterion and suggests that the Jeffreys-rule prior is a prime candidate for being a reference prior in this setting. This is the second case where asymptotic optimization of the criterion yields a prior that fits an earlier conception of noninformativity.

Theorem 2.16. *Under Conditions 0, 1, 2 and 3, for any compact set $K \subset \Theta$, we have*

$$\lim_{n \rightarrow +\infty} \left[\sup_{\Pi \in \mathcal{P}(K)} \int_K d\Pi(\theta) D(P_\theta^{\otimes n} || P_{\Pi, n}) - \frac{d}{2} \log \left(\frac{n}{2\pi e} \right) \right] = \int_K |\mathcal{I}(\theta)|^{\frac{1}{2}} d\theta. \quad (2.54)$$

Moreover, if Π^* is the Jeffreys-rule prior renormalized so that $\Pi^*(K) = 1$, then

$$\lim_{n \rightarrow +\infty} \left[\int_K d\Pi^*(\theta) D(P_\theta^{\otimes n} || P_{\Pi, n}) - \frac{d}{2} \log \left(\frac{n}{2\pi e} \right) \right] = \int_K |\mathcal{I}(\theta)|^{\frac{1}{2}} d\theta. \quad (2.55)$$

Some additional work is needed to show that the Jeffreys-rule prior is, indeed, a reference prior. For now, let us detail Clarke and Barron [1994]'s heuristic arguments supporting the above theorem. For a proof, please refer to their paper.

Heuristic arguments in support of Theorem 2.16

To simplify matters, assume $\Theta = \mathbb{R}^d$. For all θ and $\theta' \in \mathbb{R}^d$, and for all $y_1, \dots, y_n \in \mathcal{Y}$, there exists a point θ_n on the segment $[\theta, \theta']$ such that

$$\log \left(\frac{\prod_{i=1}^n p_{\theta'}(y_i)}{\prod_{i=1}^n p_\theta(y_i)} \right) = (\theta' - \theta)^\top \sum_{i=1}^n \nabla_\theta \log p_\theta(y_i) - \frac{1}{2} n (\theta' - \theta)^\top I_n(\theta_n) (\theta' - \theta), \quad (2.56)$$

where for every θ'' , $I_n(\theta'')$ is the matrix with (j, k) -th entry $\frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta_j \partial \theta_k} \log(p_{\theta''}(y_i))$. The above is a Taylor series expansion of $\log(\prod_{i=1}^n p_\theta(y_i))$.

Now, assume that the y_1, \dots, y_n are independent realizations of P_θ .

Then, for all θ'' , $I_n(\theta'')$ converges almost surely as $n \rightarrow \infty$ to the matrix with (j, k) -th entry $-\int \frac{\partial^2}{\partial \theta_j \partial \theta_k} \log p_{\theta''}(y) p_\theta(y) d\mu(y)$. If θ'' is in the neighborhood of θ , provided continuity theorems are usable, this is an approximation of the quantity $-\int \frac{\partial^2}{\partial \theta_j \partial \theta_k} \log p_\theta(y) p_\theta(y) d\mu(y)$, which is the (j, k) -th entry of the Fisher information matrix.

Let us neglect the error resulting from replacing $I_n(\theta_n)$ by $\mathcal{I}(\theta)$, the Fisher information matrix at θ .

Now let Π be a prior distribution. Let π be its Radon-Nykodym derivative with respect to the Lebesgue measure. Let us assume that π is continuous. Let $p_{\Pi}^{\otimes n}$ be defined by

$$p_{\Pi}(y_1, \dots, y_n) = \int_{\mathbb{R}^d} \prod_{i=1}^n p_{\theta'}(y_i) \pi(\theta') d\theta'. \quad (2.57)$$

p_{Π} is the Radon-Nykodym derivative with respect to μ of P_{Π} .

$$\begin{aligned} & \frac{p_{\Pi}(y_1, \dots, y_n)}{\prod_{i=1}^n p_{\theta}(y_i)} \\ &= \int_{\mathbb{R}^d} \frac{\prod_{i=1}^n p_{\theta'}(y_i)}{\prod_{i=1}^n p_{\theta}(y_i)} \pi(\theta') d\theta' \\ &= \int_{\mathbb{R}^d} \exp\left((\theta' - \theta)^{\top} \sum_{i=1}^n \nabla_{\theta} \log p_{\theta}(y_i)\right) \exp\left(-\frac{1}{2} n(\theta' - \theta)^{\top} I_n(\theta_n)(\theta' - \theta)\right) \pi(\theta') d\theta' \\ &\approx \int_{\mathbb{R}^d} \exp\left((\theta' - \theta)^{\top} \sum_{i=1}^n \nabla_{\theta} \log p_{\theta}(y_i)\right) \exp\left(-\frac{1}{2} n(\theta' - \theta)^{\top} \mathcal{I}(\theta)(\theta' - \theta)\right) \pi(\theta') d\theta'. \end{aligned} \quad (2.58)$$

Now, $(2\pi)^{-d/2} |n\mathcal{I}(\theta)|^{1/2} \exp(-\frac{1}{2} n(\theta' - \theta)^{\top} \mathcal{I}(\theta)(\theta' - \theta))$ is the density of the Normal distribution with mean θ and variance $n^{-1} \mathcal{I}(\theta)^{-1}$.

Let us define $S_n(\theta) = 1/\sqrt{n} \sum_{i=1}^n \nabla_{\theta} \log p_{\theta}(y_i)$. By the Central Limit theorem, if y_1, \dots, y_n are sampled according to P_{θ} , then the distribution of $S_n(\theta)$ is asymptotically Normal with mean $\mathbf{0}$ and variance $\mathcal{I}(\theta)$. So when $n \rightarrow \infty$, (2.58) can be approximated by averaging $\exp(\sqrt{n}(\theta' - \theta)^{\top} S_n(\theta)) \pi(\theta')$ over smaller and smaller neighborhoods of θ' . Due to π being continuous, we approximate further:

$$\begin{aligned} & \int_{\mathbb{R}^d} \exp(\sqrt{n}(\theta' - \theta)^{\top} S_n(\theta)) \exp\left(-\frac{1}{2} n(\theta' - \theta)^{\top} \mathcal{I}(\theta)(\theta' - \theta)\right) \pi(\theta') d\theta' \\ &\approx \pi(\theta) \int_{\mathbb{R}^d} \exp(\sqrt{n}(\theta' - \theta)^{\top} S_n(\theta)) \exp\left(-\frac{1}{2} n(\theta' - \theta)^{\top} \mathcal{I}(\theta)(\theta' - \theta)\right) d\theta'. \end{aligned} \quad (2.59)$$

Gathering all this,

$$\frac{p_{\Pi}(y_1, \dots, y_n)}{\prod_{i=1}^n p_{\theta}(y_i)} \approx \pi(\theta) (2\pi)^{d/2} |n\mathcal{I}(\theta)|^{-1/2} \exp\left(\frac{1}{2} S_n(\theta)^{\top} \mathcal{I}(\theta)^{-1} S_n(\theta)\right), \quad (2.60)$$

so

$$\log\left(\frac{\prod_{i=1}^n p_{\theta}(y_i)}{p_{\Pi}(y_1, \dots, y_n)}\right) \approx -\log(\pi(\theta)) + \frac{d}{2} \log\left(\frac{n}{2\pi}\right) + \frac{1}{2} \log|\mathcal{I}(\theta)| - \frac{1}{2} S_n(\theta)^{\top} \mathcal{I}(\theta)^{-1} S_n(\theta). \quad (2.61)$$

Given the distribution of $S_n(\theta)$ is asymptotically Normal with mean $\mathbf{0}$ and variance $\mathcal{I}(\theta)$ (provided y_1, \dots, y_n are independently sampled from P_{θ}), we have

$$\underbrace{\int_{\mathcal{Y}} \dots \int_{\mathcal{Y}}}_{n \text{ } f_{\mathcal{Y}}} S_n(\theta)^{\top} \mathcal{I}(\theta)^{-1} S_n(\theta) \prod_{i=1}^n p_{\theta}(y_i) d\mu(y_i) = d. \quad (2.62)$$

Therefore,

$$\begin{aligned}
D(P_\theta^{\otimes n} \| P_{\Pi, n}) &= \int_{\mathcal{Y}} \dots \int_{\mathcal{Y}} \log \left(\frac{\prod_{i=1}^n p_\theta(y_i)}{p_\Pi(y_1, \dots, y_n)} \right) \prod_{i=1}^n p_\theta(y_i) d\mu(y_i) \\
&\quad \underbrace{\hspace{10em}}_{n \int_{\mathcal{Y}}} \\
&\approx -\log(\pi(\theta)) + \frac{d}{2} \log \left(\frac{n}{2\pi e} \right) + \frac{1}{2} \log |\mathcal{I}(\theta)|. \tag{2.63}
\end{aligned}$$

Define $c = \int_{\mathbb{R}^d} |\mathcal{I}(\theta)|^{\frac{1}{2}} d\theta$ and assume c is finite. Then

$$\int_{\Theta} D(P_\theta^{\otimes n} \| P_{\Pi, n}) \pi(\theta) d\theta \approx \frac{d}{2} \log \left(\frac{n}{2\pi e} \right) + \log(c) - \int_{\Theta} \log \left(\frac{\pi(\theta)}{|\mathcal{I}(\theta)|^{\frac{1}{2}}/c} \right) \pi(\theta) d\theta. \tag{2.64}$$

Remark. The Jeffreys-rule prior can be proper or not, and is in fact usually improper. Here, though, by requiring c to be finite, we are assuming it to be proper.

The last term is the opposite of the Kullback-Leibler divergence between Π and the Jeffreys-rule prior, which is proper. To maximize this quantity, Π must be the Jeffreys-rule prior.

Equation (2.64) suggests that all priors Π that have continuous density π with respect to the Lebesgue measure make $\int_{\Theta} D(P_\theta^{\otimes n} \| P_{\Pi, n}) \pi(\theta) d\theta$ go to infinity as $n \rightarrow +\infty$. Notice that the speed of increase is the same for all continuous priors: $\log(n)$.

The Jeffreys-rule prior as reference prior

Although the Jeffreys-rule prior being a reference prior is well-known [Bernardo, 2005, Berger et al., 2009], because of our unusual definition of a reference prior, we need to make sure that the Jeffreys-rule prior satisfies this definition. The discussion in this part of the section is therefore, to the best of our knowledge, original.

Lemma 2.17 (Corollary of Clarke and Barron's theorem). *Under Conditions 0, 1, 2 and 3, let U be an open set and K a compact set such that $U \subset K \subset \Theta \subset \mathbb{R}^d$. Then for any proper prior Π with support on U ,*

$$\limsup_{n \rightarrow +\infty} \left[\int_U d\Pi(\theta) D(P_\theta^{\otimes n} \| P_{\Pi, n}) - \frac{d}{2} \log \left(\frac{n}{2\pi e} \right) \right] \leq \int_U |\mathcal{I}(\theta)|^{\frac{1}{2}} d\theta < +\infty. \tag{2.65}$$

Moreover, if Π^* is the renormalized restriction of the Jeffreys-rule prior to U , then

$$\lim_{n \rightarrow +\infty} \left[\int_U d\Pi^*(\theta) D(P_\theta^{\otimes n} \| P_{\Pi, n}) - \frac{d}{2} \log \left(\frac{n}{2\pi e} \right) \right] = \int_U |\mathcal{I}(\theta)|^{\frac{1}{2}} d\theta < +\infty. \tag{2.66}$$

Proof. Let Π be a proper prior with support on U .

If there exists a compact set $K_U \subset U$ such that $\Pi(K_U) = 1$, then Theorem 2.16 is applicable:

$$\begin{aligned}
&\limsup_{n \rightarrow +\infty} \left[\int_U d\Pi(\theta) D(P_\theta^{\otimes n} \| P_{\Pi, n}) - \frac{d}{2} \log \left(\frac{n}{2\pi e} \right) \right] \\
&= \limsup_{n \rightarrow +\infty} \left[\int_{K_U} d\Pi(\theta) D(P_\theta^{\otimes n} \| P_{\Pi, n}) - \frac{d}{2} \log \left(\frac{n}{2\pi e} \right) \right] \\
&\leq \int_{K_U} |\mathcal{I}(\theta)|^{\frac{1}{2}} d\theta \leq \int_U |\mathcal{I}(\theta)|^{\frac{1}{2}} d\theta. \tag{2.67}
\end{aligned}$$

Else, let K_0 be a compact subset of U such that $\Pi(K_0) > 0$.

Let Π_0 be the renormalized restriction of Π to K_0 (i.e., for any measurable subset A of K_0 , $\Pi_0(A) = \Pi(A)/\Pi(K_0)$) and let Π_1 be the renormalized restriction of Π to $U \setminus K_0$.

Lemma 2.6 yields that for every positive integer n ,

$$\begin{aligned} & \int_U d\Pi(\theta) D(P_\theta^{\otimes n} \| P_{\Pi, n}) \\ &= \Pi(K_0) \int_{K_0} d\Pi_0(\theta) D(P_\theta^{\otimes n} \| P_{\Pi_0, n}) + (1 - \Pi(K_0)) \int_{U \setminus K_0} d\Pi_1(\theta) D(P_\theta^{\otimes n} \| P_{\Pi_1, n}) \\ & \quad + \Pi(K_0) D(P_{\Pi_0, n} \| P_{\Pi, n}) + (1 - \Pi(K_0)) D(P_{\Pi_1, n} \| P_{\Pi, n}). \end{aligned} \quad (2.68)$$

By Lemma 2.5,

$$\begin{aligned} & \int_U d\Pi(\theta) D(P_\theta^{\otimes n} \| P_{\Pi, n}) \\ & \leq \Pi(K_0) \int_{K_0} d\Pi_0(\theta) D(P_\theta^{\otimes n} \| P_{\Pi_0, n}) + (1 - \Pi(K_0)) \int_{U \setminus K_0} d\Pi_1(\theta) D(P_\theta^{\otimes n} \| P_{\Pi_1, n}) \\ & \quad - \Pi(K_0) \log(\Pi(K_0)) - (1 - \Pi(K_0)) \log(1 - \Pi(K_0)). \end{aligned} \quad (2.69)$$

Therefore, defining $f_\Pi(K_0) := \Pi(K_0) \log(\Pi(K_0)) + (1 - \Pi(K_0)) \log(1 - \Pi(K_0))$,

$$\begin{aligned} & \int_U d\Pi(\theta) D(P_\theta^{\otimes n} \| P_{\Pi, n}) + f_\Pi(K_0) \\ & \leq \Pi(K_0) \int_{K_0} d\Pi_0(\theta) D(P_\theta^{\otimes n} \| P_{\Pi_0, n}) + (1 - \Pi(K_0)) \sup_{\Pi' \in \mathcal{P}(K)} \int_K d\Pi'(\theta) D(P_\theta^{\otimes n} \| P_{\Pi', n}) \end{aligned} \quad (2.70)$$

From this we obtain

$$\begin{aligned} & \int_U d\Pi(\theta) D(P_\theta^{\otimes n} \| P_{\Pi, n}) - \frac{d}{2} \log\left(\frac{n}{2\pi e}\right) + f_\Pi(K_0) \\ & \leq \Pi(K_0) \left[\int_{K_0} d\Pi_0(\theta) D(P_\theta^{\otimes n} \| P_{\Pi_0, n}) - \frac{d}{2} \log\left(\frac{n}{2\pi e}\right) \right] \\ & \quad + (1 - \Pi(K_0)) \left[\sup_{\Pi' \in \mathcal{P}(K)} \int_K d\Pi'(\theta) D(P_\theta^{\otimes n} \| P_{\Pi', n}) - \frac{d}{2} \log\left(\frac{n}{2\pi e}\right) \right] \end{aligned} \quad (2.71)$$

Theorem 2.16 implies that

$$\begin{aligned} & \limsup_{n \rightarrow +\infty} \left[\int_U d\Pi(\theta) D(P_\theta^{\otimes n} \| P_{\Pi, n}) - \frac{d}{2} \log\left(\frac{n}{2\pi e}\right) \right] + f_\Pi(K_0) \\ & \leq \Pi(K_0) \limsup_{n \rightarrow +\infty} \left[\int_{K_0} d\Pi_0(\theta) D(P_\theta^{\otimes n} \| P_{\Pi_0, n}) - \frac{d}{2} \log\left(\frac{n}{2\pi e}\right) \right] \\ & \quad + (1 - \Pi(K_0)) \lim_{n \rightarrow +\infty} \left[\sup_{\Pi' \in \mathcal{P}(K)} \int_K d\Pi'(\theta) D(P_\theta^{\otimes n} \| P_{\Pi', n}) - \frac{d}{2} \log\left(\frac{n}{2\pi e}\right) \right] \\ & \leq \Pi(K_0) \int_{K_0} |\mathcal{I}(\theta)|^{\frac{1}{2}} d\theta + (1 - \Pi(K_0)) \int_K |\mathcal{I}(\theta)|^{\frac{1}{2}} d\theta. \end{aligned} \quad (2.72)$$

Now, let $(K_m)_{m \in \mathbb{Z}_+}$ be an increasing sequence of compact subsets of U with limit U .

For any positive integer m ,

$$\begin{aligned} & \limsup_{n \rightarrow +\infty} \left[\int_U d\Pi(\theta) D(P_\theta^{\otimes n} \| P_{\Pi, n}) - \frac{d}{2} \log \left(\frac{n}{2\pi e} \right) \right] + f_\Pi(K_m) \\ & \leq \Pi(K_m) \int_{K_m} |\mathcal{I}(\theta)|^{\frac{1}{2}} d\theta + (1 - \Pi(K_m)) \int_K |\mathcal{I}(\theta)|^{\frac{1}{2}} d\theta. \end{aligned} \quad (2.73)$$

Because $\lim_{m \rightarrow +\infty} \Pi(K_m) = 1$, we have $\lim_{m \rightarrow +\infty} f_\Pi(K_m) = 0$ and thus

$$\limsup_{n \rightarrow +\infty} \left[\int_U d\Pi(\theta) D(P_\theta^{\otimes n} \| P_{\Pi, n}) - \frac{d}{2} \log \left(\frac{n}{2\pi e} \right) \right] \leq \lim_{m \rightarrow +\infty} \int_{K_m} |\mathcal{I}(\theta)|^{\frac{1}{2}} d\theta = \int_U |\mathcal{I}(\theta)|^{\frac{1}{2}} d\theta. \quad (2.74)$$

So the first part of the lemma is proved.

Now, let Π^* be the Jeffreys-rule prior normalized so that $\Pi^*(U) = 1$. If there exists a compact subset K_U^* of U such that $\Pi^*(K_U^*) = 1$, then by Theorem 2.16

$$\begin{aligned} & \lim_{n \rightarrow +\infty} \left[\int_U d\Pi^*(\theta) D(P_\theta^{\otimes n} \| P_{\Pi^*, n}) - \frac{d}{2} \log \left(\frac{n}{2\pi e} \right) \right] \\ & = \lim_{n \rightarrow +\infty} \left[\int_{K_U^*} d\Pi^*(\theta) D(P_\theta^{\otimes n} \| P_{\Pi, n}) - \frac{d}{2} \log \left(\frac{n}{2\pi e} \right) \right] \\ & = \int_{K_U^*} |\mathcal{I}(\theta)|^{\frac{1}{2}} d\theta = \int_U |\mathcal{I}(\theta)|^{\frac{1}{2}} d\theta. \end{aligned} \quad (2.75)$$

Else, for any compact subset $K_0 \subset U$ such that $\Pi^*(K_0) > 0$, let Π_0^* be the renormalized restriction of Π^* to K_0 Equation (2.68) implies

$$\int_U d\Pi^*(\theta) D(P_\theta^{\otimes n} \| P_{\Pi^*, n}) \geq \Pi^*(K_0) \int_{K_0} d\Pi_0^*(\theta) D(P_\theta^{\otimes n} \| P_{\Pi_0^*, n}). \quad (2.76)$$

This yields

$$\begin{aligned} & \liminf_{n \rightarrow +\infty} \left[\int_U d\Pi^*(\theta) D(P_\theta^{\otimes n} \| P_{\Pi^*, n}) - \frac{d}{2} \log \left(\frac{n}{2\pi e} \right) \right] \\ & \geq \Pi^*(K_0) \lim_{n \rightarrow +\infty} \left[\int_{K_0} d\Pi_0^*(\theta) D(P_\theta^{\otimes n} \| P_{\Pi_0^*, n}) - \frac{d}{2} \log \left(\frac{n}{2\pi e} \right) \right] \\ & = \Pi^*(K_0) \int_{K_0} |\mathcal{I}(\theta)|^{\frac{1}{2}} d\theta. \end{aligned} \quad (2.77)$$

And from there

$$\begin{aligned} \liminf_{n \rightarrow +\infty} \left[\int_U d\Pi^*(\theta) D(P_\theta^{\otimes n} \| P_{\Pi^*, n}) - \frac{d}{2} \log \left(\frac{n}{2\pi e} \right) \right] & \geq \lim_{m \rightarrow +\infty} \Pi^*(K_m) \int_{K_m} |\mathcal{I}(\theta)|^{\frac{1}{2}} d\theta \\ & = \int_U |\mathcal{I}(\theta)|^{\frac{1}{2}} d\theta. \end{aligned} \quad (2.78)$$

So the liminf is greater or equal to the limsup: this means that liminf and limsup are equal and yields the second part of the lemma.

□

Theorem 2.18. *Under Conditions 0, 1, 2 and 3, the Jeffreys-rule prior is a reference prior in the sense of Definition 2.13. If it is proper, then it is unique. If it is improper, then it is unique up to a multiplicative constant.*

Proof. Let $(U_m)_{m \in \mathbb{N}}$ be an increasing sequence of bounded open subsets of Θ with limit Θ . Therefore, for any nonnegative integer m , the closure of U_m is a compact set. This makes Lemma 2.17 applicable and proves that the Jeffreys-rule prior is a reference prior in the sense of Definition 2.13.

Both uniqueness results follow from the existence of a consistent estimator (the Maximum Likelihood estimator), which implies that Assumptions 1 and 2 hold and make Proposition 2.14 applicable. \square

2.5 Properties of reference priors

Reference priors have several interesting properties which cement their status as sensible default priors [Bernardo, 2005, Berger et al., 2009].

Theorem 2.19 (Independence from sample size). *A reference prior for a model $(P_\theta)_{\theta \in \Theta}$ remains the same for the model $(P_\theta^{\otimes n})_{\theta \in \Theta}$ regardless of $n \in \mathbb{Z}_+$.*

Proof. This is because the definition of a reference prior relies on asymptotics. \square

Theorem 2.20 (Compatibility with sufficient statistics). *If the observed data are restricted to a sufficient statistic, reference priors remain unchanged.*

Proof. Let S be a sufficient statistic. This means that there exists a family of probability distributions $(Q_s)_{s \in S(\mathcal{Y})}$ on $(\mathcal{Y}, \mathbb{Y})$ such that for all $\theta \in \Theta$, $dP_\theta(y) = d(P_\theta * S)(S(y))dQ_{S(y)}(y)$, where $P_\theta * S$ is the push-forward measure of P_θ by S (if Y is a random variable that follows the distribution P_θ , then $S(Y)$ follows $P_\theta * S$).

For every positive integer n ,

$$\begin{aligned}
D(P_\theta^{\otimes n} || P_{\Pi_n}) &= \int_{\mathcal{Y}^n} dP_\theta^{\otimes n}(y_1, \dots, y_n) \log \left(\frac{\prod_{i=1}^n dP_\theta(y_i)}{\int_{\Theta} d\Pi(\theta) \prod_{i=1}^n dP_\theta(y_i)} \right) \\
&= \int_{\mathcal{Y}^n} dP_\theta^{\otimes n}(y_1, \dots, y_n) \log \left(\frac{\prod_{i=1}^n d(P_\theta * S)(S(y_i))dQ_{S(y_i)}(y_i)}{\int_{\Theta} d\Pi(\theta) \prod_{i=1}^n d(P_\theta * S)(S(y_i))dQ_{S(y_i)}(y_i)} \right) \\
&= \int_{\mathcal{Y}^n} dP_\theta^{\otimes n}(y_1, \dots, y_n) \log \left(\frac{d(P_\theta * S)^{\otimes n}(S(y_1), \dots, S(y_n))}{\int_{\Theta} d\Pi(\theta) d(P_\theta * S)^{\otimes n}(S(y_1), \dots, S(y_n))} \right) \\
&= \int_{S(\mathcal{Y})^n} d(P_\theta * S)^{\otimes n}(s_1, \dots, s_n) \log \left(\frac{d(P_\theta * S)^{\otimes n}(s_1, \dots, s_n)}{\int_{\Theta} d\Pi(\theta) d(P_\theta * S)^{\otimes n}(s_1, \dots, s_n)} \right).
\end{aligned} \tag{2.79}$$

\square

Theorem 2.21 (Consistency under reparametrization). *If Θ' is a metric space and $f : \Theta \rightarrow \Theta'$ is a measurable bijection whose inverse f^{-1} is also measurable, and if Π^* is a reference prior for the model $(P_\theta)_{\theta \in \Theta}$, then the push-forward measure $\Pi^* * f$ is a reference prior for the model $(P_{f^{-1}(\theta')})_{\theta' \in \Theta'}$.*

Proof. This follows from the change-of-variable formula. For any measurable function $g : \Theta \rightarrow [0, +\infty)$,

$$\int_{\Theta} d\Pi(\theta)g(\theta) = \int_{\Theta'} d(\Pi * f)(\theta') g(f^{-1}(\theta')). \quad (2.80)$$

In particular, for any positive integer n and any measurable set $A \in \mathbb{Y}^{\otimes n}$,

$$P_{\Pi, n}(A) = \int_{\Theta} d\Pi(\theta)P_{\theta}^{\otimes n}(A) = \int_{\Theta'} d(\Pi * f)(\theta')P_{f^{-1}(\theta')}^{\otimes n}(A). \quad (2.81)$$

□

The latter two properties can be interpreted this way: if it exists, a reference prior only cares about the model as a family of probability distributions. The parameter is only an index with no intrinsic meaning (consistency under reparametrization). What matters is how the probability distributions in the model differ from each other and any common parts are disregarded (compatibility with sufficient statistics).

2.6 Examples

The first two examples – location and scale models – are drawn from Bernardo [2005] and Berger et al. [2009].

Proposition 2.22 (Location models). *If $\Theta = \mathcal{Y} = \mathbb{R}^d$ ($d \in \mathbb{Z}_+$), with \mathbb{R}^d endowed with the Borel σ -algebra $\mathcal{B}(\mathbb{R}^d)$, if for any $\theta, \theta' \in \mathbb{R}^d$ and any measurable subset A of \mathbb{R}^d , $P_{\theta'}(A) = P_{\theta}(A + \theta - \theta')$, and if Assumption 2 holds, then any reference prior is proportional to the Lebesgue measure.*

Proof. Let Π^* be a reference prior for such a model. For any $\theta_0 \in \mathbb{R}^d$, define the statistic $S_{\theta_0} : \mathbb{R}^d \rightarrow \mathbb{R}^d; y \mapsto y + \theta_0$. S_{θ_0} is a sufficient statistic, so compatibility with sufficient statistics implies Π^* stays reference prior for the model $(P_{\theta} * S_{\theta_0})_{\theta \in \mathbb{R}^d}$. However, for all $\theta \in \mathbb{R}^d$, $P_{\theta} * S_{\theta_0} = P_{\theta + \theta_0}$. In other words, S defines a reparametrization. Therefore consistency under reparametrization implies the reference prior for $(P_{\theta + \theta_0})_{\theta \in \mathbb{R}^d}$ corresponding to Π^* is the prior $\Pi_{\theta_0}^*$ defined for every measurable set A by $\Pi_{\theta_0}^*(A) = \Pi^*(A - \theta_0)$.

Gathering this, we obtain $\Pi_{\theta_0}^* = \Pi^*$, and this holds for any $\theta_0 \in \mathbb{R}^d$. This implies that Π^* is proportional to the Lebesgue measure.

□

Proposition 2.23 (Scale models). *Assume $\mathcal{Y} = \mathbb{R}^d$ ($d \in \mathbb{Z}_+$) and $\Theta = (0, +\infty)$ are both endowed with the Euclidean or any equivalent distance. If for any $\theta, \theta' \in (0, +\infty)$ and any measurable subset A of \mathbb{R}^d , $P_{\theta'}(A) = P_{\theta}(\frac{\theta'}{\theta}A)$, and if Assumption 2 holds, then any reference prior is proportional to $\theta^{-1}d\theta$.*

Proof. In such a model, there must exist a sufficient statistic with value on $(0, +\infty)$ for which θ works as a scale parameter (to build it, first construct a statistic whose distribution under any value of θ is isotropic and then take its Euclidean norm). Reference priors are compatible with sufficient statistics, so let us replace the model with this one. This allows us to apply the logarithm to the observation space, mapping it to \mathbb{R} . Taking $\log(\theta)$ as the new parametrization, the model becomes a location model and any reference prior is proportional

to the Lebesgue measure. Given reference priors are consistent under reparametrization, any reference prior on θ is proportional to the push-forward of the Lebesgue measure by the exponential function, that is the measure $\theta^{-1}d\theta$. \square

Location and scale models are remarkable because invariance properties yield the reference prior. We now study an example where no such property is available. This example appears in Ren et al. [2012], but the reference prior is derived differently.

Example 3. Consider a model $(P_\theta)_{\theta \in (0, +\infty)}$ where each P_θ is absolutely continuous with respect to the Lebesgue measure on the unit sphere $S^{n-1} := \{\mathbf{y} \in \mathbb{R}^n : \mathbf{y}^\top \mathbf{y} = 1\}$. The Radon-Nykodym derivative p_θ of P_θ (i.e. the likelihood function) is given by:

$$p_\theta(\mathbf{y}) = \left(\frac{2\pi^{\frac{n}{2}}}{\Gamma(\frac{n}{2})} \right)^{-1} |\Sigma_\theta|^{-\frac{1}{2}} (\mathbf{y}^\top \Sigma_\theta^{-1} \mathbf{y})^{-\frac{n}{2}}, \quad (2.82)$$

where Σ_θ is a symmetric positive definite matrix and is twice continuously differentiable as a function of $\theta \in (0, +\infty)$. Therefore Condition 0 of the regular case holds. Further, let us assume that $\Sigma_\theta = \Sigma_{\theta'} \Rightarrow \theta = \theta'$, which means that Condition 3 holds. Since the sphere S^{n-1} is a compact set, Conditions 1 and 2 hold as well, so Theorem 2.18 is applicable. The reference prior is the Jeffreys-rule prior.

Define a matrix $\sqrt{\Sigma_\theta}$ such that $\Sigma_\theta = \sqrt{\Sigma_\theta} \sqrt{\Sigma_\theta}^\top$ (use for instance the Cholesky decomposition). Now let \mathbf{Y} be a random variable following P_θ . Then $f_\theta(\mathbf{Y}) := \sqrt{\Sigma_\theta}^{-1} \mathbf{Y} / \|\sqrt{\Sigma_\theta}^{-1} \mathbf{Y}\|$ follows the Uniform distribution on the sphere S^{n-1} .

Besides,

$$\partial_\theta (\log p_\theta(\mathbf{y})) = -\frac{n}{2} f_\theta(\mathbf{y})^\top \mathbf{M}_\theta^\Sigma f_\theta(\mathbf{y}) + C_\theta, \quad (2.83)$$

where

$$C_\theta := -\log \left(\frac{2\pi^{\frac{n}{2}}}{\Gamma(\frac{n}{2})} \right) - \frac{1}{2} \partial_\theta |\Sigma_\theta|^{-\frac{1}{2}} \quad \text{and} \quad \mathbf{M}_\theta^\Sigma := \left(\sqrt{\Sigma_\theta} \right)^\top \partial_\theta (\Sigma_\theta^{-1}) \sqrt{\Sigma_\theta}. \quad (2.84)$$

Therefore

$$\text{Var} [\partial_\theta (\log p_\theta(\mathbf{Y}))] = \frac{n^2}{4} \text{Var} \left[f_\theta(\mathbf{Y})^\top \mathbf{M}_\theta^\Sigma f_\theta(\mathbf{Y}) \right]. \quad (2.85)$$

As \mathbf{M}_θ^Σ is a symmetric matrix, the spectral theorem guarantees the existence of a diagonal matrix $\mathbf{\Lambda}_\theta^\Sigma$ and an orthogonal matrix \mathbf{O}_θ^Σ such that $\mathbf{M}_\theta^\Sigma = (\mathbf{O}_\theta^\Sigma)^\top \mathbf{\Lambda}_\theta^\Sigma (\mathbf{O}_\theta^\Sigma)$, with the diagonal coefficients of $\mathbf{\Lambda}_\theta^\Sigma$ being the eigenvalues of \mathbf{M}_θ^Σ . Setting $\mathbf{U}_0 := (\mathbf{O}_\theta^\Sigma) \mathbf{U}$, we can now compute $\text{Var}[\mathbf{U}_0^\top \mathbf{\Lambda}_\theta^\Sigma \mathbf{U}_0] = \text{Var}[\mathbf{U}^\top \mathbf{M}_\theta^\Sigma \mathbf{U}]$, \mathbf{U}_0 following the uniform distribution on S^{n-1} .

Let $(\lambda_i)_{1 \leq i \leq n}$ be the eigenvalues of \mathbf{M}_θ^Σ . We can write $\text{Var}[\mathbf{U}_0^\top \mathbf{\Lambda}_\theta^\Sigma \mathbf{U}_0] = \text{Var}[\sum_{1 \leq i \leq n} \lambda_i X_i]$, where X_i ($1 \leq i \leq n$) are nonnegative identically distributed random variables such that $\sum_{1 \leq i \leq n} X_i = 1$.

$$\text{Var} \left[\sum_{i=1}^n \lambda_i X_i \right] = \text{Var}[X_1] \sum_{i=1}^n \lambda_i^2 + 2 \text{Cov}(X_1, X_2) \sum_{1 \leq i < j \leq n} \lambda_i \lambda_j. \quad (2.86)$$

Because $\mathbb{E}[X_1] = \frac{1}{n}$, $\text{Cov}(X_1, X_2) = -1/(n-1) \text{Var}[X_1]$.

$$\begin{aligned}
\frac{\text{Var}[\sum_{i=1}^n \lambda_i X_i]}{\text{Var}[X_1]} &= \sum_{i=1}^n \lambda_i^2 - \frac{1}{n-1} \sum_{i=1}^n \lambda_i \sum_{j \neq i} \lambda_j \\
&= \left(1 + \frac{1}{n-1}\right) \left(\text{Tr} \left[(\mathbf{M}_\theta^\Sigma)^2 \right] - \frac{1}{n} \text{Tr} \left[\mathbf{M}_\theta^\Sigma \right]^2 \right) \\
&= \left(1 + \frac{1}{n-1}\right) \left(\text{Tr} \left[\left(\left(\frac{\partial}{\partial \theta} \boldsymbol{\Sigma}_\theta \right) \boldsymbol{\Sigma}_\theta^{-1} \right)^2 \right] - \frac{1}{n} \text{Tr} \left[\left(\frac{\partial}{\partial \theta} \boldsymbol{\Sigma}_\theta \right) \boldsymbol{\Sigma}_\theta^{-1} \right]^2 \right).
\end{aligned} \tag{2.87}$$

The reference prior therefore has density π with respect to the Lebesgue measure on $(0, +\infty)$, where π is defined by:

$$\pi(\theta) \propto \sqrt{\text{Tr} \left[\left(\left(\frac{\partial}{\partial \theta} \boldsymbol{\Sigma}_\theta \right) \boldsymbol{\Sigma}_\theta^{-1} \right)^2 \right] - \frac{1}{n} \text{Tr} \left[\left(\frac{\partial}{\partial \theta} \boldsymbol{\Sigma}_\theta \right) \boldsymbol{\Sigma}_\theta^{-1} \right]^2}. \tag{2.88}$$

Remark. Consider the case where $\theta \in (0, +\infty)^r$, r being a greater than 1 integer. To emphasize that θ is multidimensional, let us denote it by $\boldsymbol{\theta}$. The reference prior is still the Jeffreys rule prior: it is the square root of the determinant of the Fisher information matrix $\mathcal{I}(\boldsymbol{\theta})$ whose (i, j) -th element is

$$\text{Tr} \left[\left(\frac{\partial}{\partial \theta_i} (\boldsymbol{\Sigma}_\theta \boldsymbol{\Sigma}_\theta^{-1}) \right) \left(\frac{\partial}{\partial \theta_j} (\boldsymbol{\Sigma}_\theta \boldsymbol{\Sigma}_\theta^{-1}) \right) \right] - \frac{1}{n} \text{Tr} \left[\frac{\partial}{\partial \theta_i} (\boldsymbol{\Sigma}_\theta \boldsymbol{\Sigma}_\theta^{-1}) \right] \text{Tr} \left[\frac{\partial}{\partial \theta_j} (\boldsymbol{\Sigma}_\theta \boldsymbol{\Sigma}_\theta^{-1}) \right]. \tag{2.89}$$

This results from polarization formula $\text{Cov}(A, B) = 1/4(\text{Var}(A + B) - \text{Var}(A - B))$.

Let us now consider a particular example that is a hybrid between location and scale models. It is often referenced in the literature. See for instance Robert [2007, section 3.5].

Example 4. Take $\theta = (\beta, \sigma)$ and let $P_{\beta, \sigma} = \mathcal{N}(\beta, \sigma^2)$ be the one-dimensional Normal model. This is the sort of regular case where the reference prior is the Jeffreys-rule prior. The Fisher information matrix is diagonal with both diagonal elements equal to σ^{-2} . The Jeffreys-rule prior has therefore density σ^{-2} . Let us examine this a little further. If σ were known, then this would be a location model and the reference prior on β would have density $\propto 1$. If β were known, then this would be a scalar model and the reference prior would have density $\propto \sigma^{-1}$.

Consider n independent real-valued observations y_1, \dots, y_n of the model.

If β were known but not σ , then the posterior distribution of $\sum_{y=1}^n (y_i - \beta)^2 / \sigma^2$ (i.e. its distribution knowing all y_i ($1 \leq i \leq n$) and β) would be the chi-squared distribution with n degrees of freedom.

In our case, both β and σ^2 are unknown. Let us define $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$. Then the posterior distribution of $\sum_{y=1}^n (y_i - \bar{y})^2 / \sigma^2$ (i.e. its distribution knowing all y_i ($1 \leq i \leq n$) but not β) is also the chi-squared distribution with n degrees of freedom. In other words, using the reference prior in both situations implies that by simply substituting the empirical mean \bar{y} to the actual mean β , we are able to reach the same state of knowledge about σ as if we actually knew β .

Jeffreys [1961] suggested using what is now called the independence Jeffreys prior distribution, which is the joint prior distribution on β and σ obtained by taking the product of the Jeffreys-rule (reference) prior on β when σ is known and of the Jeffreys-rule (reference) prior on σ

when β is known, thus yielding the joint prior distribution σ^{-1} . Then the posterior distribution of $\sum_{y=1}^n (y_i - \bar{y})^2 / \sigma^2$ is the chi-squared distribution with $n - 1$ degrees of freedom, which acknowledges the loss of information on σ^2 when β is unknown.

This example suggests that reference priors, as they are currently defined, are not necessarily adequate for dealing with multiparametric settings. It also suggests a method for solving the problem.

2.7 Reference priors for multiparametric models

Previously, we dealt with the case of one-parameter models. Naturally, in the case of multi-parameter models, it is always possible to view the list of parameters as one big multi-dimensional parameter, so as to reduce the problem to the one already tackled. Such a choice often leads to unfortunate results however, in the sense that the obtained reference prior seems intuitively unsatisfactory. More precisely, there exist several multi-parameter models in the literature in which the reference prior, as defined above, has undesirable statistical properties. See Berger et al. [2015] for a large review of such situations. In fact, the authors state: “ We actually know of no multivariable example in which we would recommend the Jeffreys-rule prior. In higher dimensions, the prior always seems to be either ‘too diffuse’ [...] or ‘too concentrated’ ”. And, as mentioned before, the reference prior is the Jeffreys-rule prior in regular cases.

Let $\Theta_1 \times \dots \times \Theta_r$ ($r \in \mathbb{Z}_+$) be the parametric space. The “reference prior algorithm”, which was first developed by Bernardo [1979a], first requires an ordering of the parameters. Let us consider them ordered in the following way:

$$\theta_1 \in \Theta_1 \prec \dots \prec \theta_r \in \Theta_r$$

We denote by $\boldsymbol{\theta}_{[j]}$ the collection $(\theta_j, \theta_{j+1}, \dots, \theta_r)$. Define $\mathbb{Y}_1 := \mathbb{Y}$.

Before delving into the details, let us explain the heuristic behind the algorithm. The idea is quite simple.

The reference prior is only defined in one-parameter settings. Let us therefore fix all parameters except one. Then we can compute the reference prior on the remaining parameter conditionally to all fixed parameters. Then, having defined a probability distribution on one of the parameters, we can integrate it out of the model.

We now find ourselves with a model with one less parameter. The same procedure can be applied again and again until none is left. The reference prior is then defined to be the product of the priors derived at every stage.

Proper reference prior

If possible, follow this algorithm: for every integer $j \in \llbracket 1, r \rrbracket$,

1. For all $\boldsymbol{\theta}_{[j+1]} \in \Theta_{j+1} \times \dots \times \Theta_r$, compute $\Pi_j(\cdot | \boldsymbol{\theta}_{[j+1]})$. It is defined as the reference prior on θ_j with respect to the model $(P_{\boldsymbol{\theta}_{[j]}})_{\boldsymbol{\theta}_{[j+1]} \in \Theta_{j+1} \times \dots \times \Theta_r}$ when $\boldsymbol{\theta}_{[j+1]}$ is assumed known. If, for some $\boldsymbol{\theta}_{[j+1]} \in \Theta_{j+1} \times \dots \times \Theta_r$, it does not exist, abort.
2. If there exists $\boldsymbol{\theta}_{[j+1]} \in \Theta_{j+1} \times \dots \times \Theta_r$ such that $\Pi_j(\cdot | \boldsymbol{\theta}_{[j+1]})$ is improper, abort.

3. For all $\boldsymbol{\theta}_{[j+1]} \in \Theta_{j+1} \times \dots \times \Theta_r$, compute $P_{\boldsymbol{\theta}_{[j+1]}} := \int_{\Theta_j} P_{\boldsymbol{\theta}_{[j]}} d\Pi_j(\theta_j | \boldsymbol{\theta}_{[j+1]})$.

If Assumption 1 holds at every stage, no choice is required from the user aside from the initial ordering of the parameters.

The reference prior based on the chosen ordering is then defined as the product $d\Pi(\boldsymbol{\theta}_{[1]}) = \prod_{j=1}^r d\Pi_j(\theta_j | \boldsymbol{\theta}_{[j+1]})$.

The general idea behind this algorithm is to emulate the success of the one-dimensional reference prior by creating such a situation for the parameter that is of “greatest interest”. To achieve this, “nuisance” parameters are integrated out of the model. To simplify, consider the 2-parameter case, in which there is only one nuisance parameter. In order to create a one-parameter model for the parameter of interest, the nuisance parameter must be integrated out of the original model. But to achieve this, a prior distribution on the nuisance parameter conditional to the parameter of interest must be defined. Hence the idea to define it as the reference prior for the model where the parameter of interest is known.

It must be stressed that this reference prior for multiparametric models is not a further extension of the concept of reference prior like Definition 2.13 was to Definition 2.4 and Definition 2.4 to Definition 2.2. It is an altogether different notion because it introduces structure in the model in the form of parameter ordering. Yang and Berger [1996] list different reference priors obtained with different parameter orderings for a large number of statistical models.

Applying this “reference prior algorithm” results in a “marginal reference prior” on the parameter of interest that is in fact the reference prior (in the sense of Definition 2.13) of a specific mixture model [Marin et al., 2005, Mengersen et al., 2011]. So the reference prior resulting from application of the reference prior algorithm is the solution of a very different optimization problem than the one that is solved by the reference prior in the sense of Definition 2.13.

Notice that the algorithm contains two abortive checks per step. The first is intrinsically linked to reference prior theory, which does not provide guarantees of existence. The second is due to the fact that a valid statistical model can only contain probability distributions: it does not tolerate infinite measures! In other words, a prior can well be improper, but a likelihood function cannot.

Note however that aborting at stage 2 is not a problem if $j = r$, because the last step has been reached. So all “conditional” reference priors must be proper, but the last, “marginal” one can be improper. As will be seen in the next subsection, this propriety requirement on the conditionals can be removed, making stage 2 of the algorithm irrelevant.

Improper reference prior

The algorithm above can be altered to remove one of the two abortive checks. Changes are written in **bold font**. First, **define** $\mathbb{Y}_1 := \mathbb{Y}$. Then, for every integer $j \in \llbracket 1, r \rrbracket$,

1. For all $\boldsymbol{\theta}_{[j+1]} \in \Theta_{j+1} \times \dots \times \Theta_r$, compute $\Pi_j(\cdot | \boldsymbol{\theta}_{[j+1]})$. It is defined as the reference prior on θ_j with respect to the model $(P_{\boldsymbol{\theta}_{[j]}})_{\boldsymbol{\theta}_{[j+1]} \in \Theta_{j+1} \times \dots \times \Theta_r}$ when $\boldsymbol{\theta}_{[j+1]}$ is assumed known. If, for some $\boldsymbol{\theta}_{[j+1]} \in \Theta_{j+1} \times \dots \times \Theta_r$, it does not exist, abort.

2. If there exists $\theta_{[j+1]} \in \Theta_{j+1} \times \dots \times \Theta_r$ such that $\Pi_j(\cdot|\theta_{[j+1]})$ is improper, **do the following. Choose a sub- σ -algebra $\mathbb{Y}_{j+1} \subset \mathbb{Y}_j$ such that for every $A \in \mathbb{Y}_{j+1}$, $P_{\theta_{[j]}}(A)$ only depends on $\theta_{[j+1]}$. Define $P_{\theta_{[j+1]}}$ on $\mathbb{Y}_{j+1} \subset \mathbb{Y}_j$ as the restriction of $P_{\theta_{[j]}}$ to \mathbb{Y}_{j+1} .** Else, for all $\theta_{[j+1]} \in \Theta_{j+1} \times \dots \times \Theta_r$, compute $P_{\theta_{[j+1]}} := \int_{\Theta_j} P_{\theta_{[j]}} d\Pi_j(\theta_j|\theta_{[j+1]})$.

As written above, this altered algorithm – henceforth called “reference prior algorithm” – has the advantage of removing one aborting check. In cases where the first algorithm works, this one is functionally identical and produces the same result. In cases where improper priors are encountered before the last stage ($j = r$) however, it requires a choice. Therefore, even if Assumption 1 holds at every stage, the reference prior is only defined with reference to a given filtering $\mathbb{Y} = \mathbb{Y}_1 \supset \dots \supset \mathbb{Y}_r$.

Remark. The main problem of this algorithm, besides the lack of uniqueness of the reference prior, is the risk that at some stage, the only possible choice for \mathbb{Y}_j might be the trivial σ -algebra. One possible solution is to consider the model with n independent observations instead of one, in order to manipulate richer σ -algebras. An example is shown below.

The altered algorithm does not fit the approach that is most often used to deal with impropriety. The idea recommended by Berger and Bernardo [1992] is to choose for every $j \in \llbracket 1, r \rrbracket$ an increasing sequence of compact subspaces $(\Theta_j^{(k)})_{k \in \mathbb{N}}$ such that $\bigcup_{k \in \mathbb{N}} \Theta_j^{(k)} = \Theta_j$ with the understanding that the reference prior on each compact is proper. This happens when the reference prior is the Jeffreys-rule prior, for example. Then, it is possible to use the first algorithm on each $\Theta_1^{(k)} \times \dots \times \Theta_r^{(k)}$, compute the associated reference posterior, take the limit of the posterior when $k \rightarrow +\infty$ and then define the reference prior as the prior which, using Bayes’ rule, would yield the reference posterior. This increasing-compact-sequence approach is what Berger et al. [2001] call *asymptotic marginalization*, whereas the independent-sub- σ -algebra approach detailed above is a formalization of the *exact marginalization* they recommend for Kriging models.

Unfortunately, asymptotic marginalization does not guarantee uniqueness of the reference prior (even for a given parameter ordering) any more than exact marginalization does. The choice of sequences of compact subspaces can influence the result. Because this thesis deals with Kriging and Kriging-derived models, all reference priors discussed in the dissertation are obtained through exact marginalization.

Example 5. *Let us come back to the earlier Normal $\mathcal{N}(\beta, \sigma^2)$ example. To apply the reference prior algorithm, we need to define an ordering on the parameters β and σ . Let us choose the ordering $\beta \prec \sigma$.*

At the first stage of the algorithm, the reference prior on β (knowing σ) is, up to a multiplicative constant, the Lebesgue measure (cf. earlier example). This is an improper prior, so we need to find a sub- σ -algebra of $\mathcal{B}(\mathbb{R})$ on which the model does not depend on β . To do this, let us consider the model with $n \geq 2$ observations. We are looking for a sub- σ -algebra of $\mathcal{B}(\mathbb{R})^{\otimes n} = \mathcal{B}(\mathbb{R}^n)$.

Let us consider the statistic $S_n : \mathbb{R}^n \rightarrow \mathbb{R}^{n-1}$ defined by $S_n((y_1, \dots, y_n)^\top) = (y_2 - y_1, \dots, y_n - y_1)^\top$. If \mathbf{Y} is a random variable following $\mathcal{N}(\beta, \sigma^2)^{\otimes n}$, $S_n(\mathbf{Y})$ follows the $(n - 1)$ -variate Normal distribution $\mathcal{N}(\mathbf{0}_{n-1}, \sigma^2(\mathbf{1}\mathbf{1}^\top + \mathbf{I}_{n-1}))$, with $\mathbf{1} = (1, 1, \dots, 1)^\top$ and with \mathbf{I}_{n-1} being the identity matrix. Therefore the sub- σ -algebra of $\mathcal{B}(\mathbb{R}^n)$ spanned by S_n does not depend

on β . Moreover, the restriction of the model to this sub- σ -algebra is a simple scale model. Therefore the reference prior on σ has density proportional $1/\sigma$ with respect to the Lebesgue measure on $(0, +\infty)$.

Gathering this, the reference prior associated to the decomposition we chose is $1/\sigma d\beta d\sigma$. Interestingly, this is the “independence Jeffreys prior” that was recommended by Jeffreys [1961] for this situation.

Finally, note that the same reference prior could have been obtained by choosing the other ordering: $\sigma \prec \beta$. At the first stage of the algorithm, when β is taken to be known, the model is a scale model so the reference prior on σ knowing β has density proportional to $1/\sigma$ with respect to the Lebesgue measure. This is an improper prior, so we need to find a sub- σ -algebra on which the model does not depend on σ . The statistic $T_n : (y_1, \dots, y_n)^\top \mapsto \sqrt{n}\bar{y}_n/\hat{\sigma}_n$ (with $\bar{y}_n = n^{-1} \sum_{i=1}^n y_i$ and $\hat{\sigma}_n^2 = n^{-1} \sum_{i=1}^n (y_i - \bar{y}_n)^2$) spans such a sub- σ -algebra. If \mathbf{Y} follows $\mathcal{N}(\beta, \sigma^2)$, then $T_n(\mathbf{y})$ follows a noncentral Student t -distribution with $n-1$ degrees of freedom and noncentrality parameter β . The resulting model is a location model, so the reference prior on β is proportional to the Lebesgue measure.

This example is important, but it should be noted that cases where the same reference prior is obtained regardless of the ordering of the parameters are the exception rather than the norm. Berger et al. [2015] explores several methods to find an “overall reference priors” which would be adequate regardless of the ordering of the parameters.

Some final thoughts

In the Normal example, the Jeffreys-rule prior is considered inadequate because it has undesirable inferential properties. It is, though, the reference prior for the model if $(\beta, \sigma)^\top$ is viewed as a single multidimensional parameter. So why does the reference prior behave “badly” in this instance? First, let us observe that it does not differ from the independence Jeffreys prior regarding the mean parameter β . The difference between the two concerns the standard deviation σ . The Jeffreys-rule prior favors small standard deviations, which lead to concentrated samples. Recalling the game-theoretic interpretation from Clarke and Barron [1994], Nature is giving a hint to the Statistician – the standard deviation is likely to be small. Why is it in Nature’s interest to do so? Because this forces the Statistician to predict a concentrated sample, even though the Statistician does not know around which value it should concentrate. If Nature’s prior favored greater standard deviations, the Statistician could “spread the prediction” around and would be more likely to cover the true mean. The Jeffreys-rule prior is truly noninformative in the sense that it makes prediction most difficult *before* viewing the data, even though the Statistician would have less trouble inferring σ *after*.

The Jeffreys independence prior is better not because it is less but because it is *more* informative than the Jeffreys-rule prior. It contains information about the structure of the model by clearly separating mean from standard deviation. It turns out that being invariant under reparametrization is not necessarily an advantage, because it can take away useful information that is contained in the structure of the parametrization.

Instead of being invariant under any reparametrization, a reference prior obtained by application of the reference prior algorithm is invarian under “sensible” reparametrizations, i.e.

reparametrizations that respect the structure. Consider a reference prior for $\theta_1 \prec \dots \prec \theta_r$ (with associated parameter space $\Theta_1 \times \dots \times \Theta_r$). If for every integer $i \in \llbracket 1, r \rrbracket$, f_i is a measurable bijection $\Theta_i \rightarrow \Theta'_i$ with measurable inverse, then this reference prior is invariant under the reparametrization $(\theta_1, \dots, \theta_r) \mapsto (f_1(\theta_1), \dots, f_r(\theta_r))$.

The claim in the introduction to this chapter that a reference prior should be used when no prior information is available is therefore not wholly correct. If the user can define a parametrization, it implies that some information is available and should be used.

Deriving a reference prior is not about producing the least informative prior, but rather about controlling the information contained in the prior. From this point onward, following Berger and Bernardo, we call it “objective” rather than noninformative.

However, this word should not lead to the conclusion that a reference prior is “better” than another because of its objectivity. Depending on whether the goal is prediction or inference, and on what is being inferred, other priors may produce better results. Tuyl et al. [2016] show reference (or Jeffreys-rule) priors are not adequate for the study of rare events. Indeed, this is not what they were designed for: by definition, rare events play a small role in the optimization criterion of reference priors.

Currently, research focuses on choosing the prior in order to obtain a specific, desired, effect on the posterior. Xueou et al. [2018] propose an elicitation mechanism for the prior Π based on history matching in order to make relevant summary statistics [Marin et al., 2014] of data produced by P_Π fit an expected behavior. Computation complexity is dealt with by using Approximate Bayesian Computation [Marin et al., 2012]. Simpson et al. [2017] introduce a new framework called “Penalized Complexity” priors as an attempt to formalize the principle of Occam’s razor for prior elicitation – again with reference to P_Π ; Robert and Rousseau [2017] note that the authors attempt to circumvent the subjectivity inherent to the ordering of parameters by dividing them into independent components. The original developers of the reference prior paradigm are also concerned with this issue and attempt to address it through the notion of “overall reference priors” [Berger et al., 2015]. Another direction for research of reasonable default priors is the Robust Bayesian framework; for a recent review of this field, see Watson and Holmes [2016].

Chapter 3

Propriety of the reference posterior distribution in Gaussian Process regression

This chapter draws on the article Muré [2018a].

Abstract

In a seminal article, Berger et al. [2001] compare several objective prior distributions for the parameters of Gaussian Process regression models with isotropic correlation kernel. The reference prior distribution stands out among them insofar as it always leads to a proper posterior. They prove this result for rough correlation kernels - Spherical, Exponential with power $q < 2$, Matérn with smoothness $\nu < 1$. This chapter provides a proof for smooth correlation kernels - Exponential with power $q = 2$, Matérn with smoothness $\nu \geq 1$, Rational Quadratic.

Résumé

Dans un article fondamental, Berger et al. [2001] comparent plusieurs lois *a priori* objectives sur les paramètres de modèles de régression par processus gaussiens avec noyau de corrélation isotrope. Le prior de référence se distingue parmi eux en cela qu'il mène systématiquement à un posterior propre. Ils démontrent ce résultat pour des noyaux de corrélation rugueux - noyau sphérique, exponentiel avec puissance $q < 2$, Matérn de régularité $\nu < 1$. Ce chapitre fournit une preuve valable pour des noyaux de corrélation réguliers - exponentiel de puissance $q = 2$, Matérn de régularité $\nu \geq 1$, rationnel quadratique.

3.1 Introduction

In a very influential paper, Berger et al. [2001] pioneered the field of Objective Bayesian analysis of spatial models. Previous works [De Oliveira et al., 1997, Stein, 1999] had noted that commonly used noninformative priors sometimes failed to yield proper posteriors, but Berger et al. [2001] were the first to thoroughly investigate the issue. Among several prior distributions - truncated priors, vague priors, Jeffreys-rule and independence Jeffreys prior - they showed that the reference prior (i.e., the reference prior with the parameter ordering

they chose) is the most satisfying choice for a default prior distribution. This is due in no small part to the fact that, in the wide variety of cases studied by Berger et al. [2001], it systematically yields a proper posterior distribution. In this article, we complete their proof of this property.

Interestingly, Berger et al. [2001] found that the corresponding reference prior obtained through asymptotic marginalization does not share this property. In this chapter and in the rest of this dissertation, all reference priors discussed are obtained through exact marginalization.

Section 3.2 describes the Gaussian Process models studied by Berger et al. [2001]. Section 3.3 shows that the proof of reference posterior propriety provided by Berger et al. [2001] only applies to those with rough correlation kernels – Spherical, Exponential with power $q < 2$, Matérn with smoothness $\nu < 1$. Section 3.4 contains the core of this chapter: a proof of Theorem 3.9 which asserts that the reference prior leads to a proper posterior for models with smoother correlation kernels – Exponential with power $q = 2$, Matérn with smoothness $\nu \geq 1$, Rational Quadratic.

Because it is difficult to obtain a satisfying default prior distribution which consistently yields a proper posterior, it is important to ascertain that the reference prior actually does. Indeed, a vast literature [Paulo, 2005, Ren et al., 2012, Kazianka and Pilz, 2012, Ren et al., 2013, Gu et al., 2018] builds upon Berger et al. [2001]’s result and depends on it.

3.2 Setting

Berger et al. [2001] consider models of Gaussian Process regression, also known as Universal Kriging, with isotropic autocorrelation kernels. Because isotropy is key, define $\|\cdot\|$ as the usual Euclidean norm if applied to a vector and as the Frobenius norm if applied to a matrix. In Universal Kriging, an unknown mapping from a spatial domain $\mathcal{D} \subset \mathbb{R}^r$ ($r \in \mathbb{Z}_+$) to \mathbb{R} is assumed to be a realization of a Gaussian process Y . The mean function f of the Gaussian process is assumed to belong to some known vector space \mathcal{F}_p of dimension $p \in \mathbb{N}$. If p is non-zero, once a basis $(f_j)_{j \in \llbracket 1, p \rrbracket}$ of \mathcal{F}_p has been set, f can be parametrized by $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top \in \mathbb{R}^p$ such that $f = \sum_{j=1}^p \beta_j f_j$.

$Y - f$ is assumed in the model to be an isotropic Gaussian process based on an autocorrelation kernel K . K is a mapping $[0, +\infty) \rightarrow \mathbb{R}$ such that for any positive integer n and any collection of n distinct points $(\mathbf{x}^{(i)})_{i \in \llbracket 1, n \rrbracket}$ within \mathcal{D} , the symmetric $n \times n$ matrix $\boldsymbol{\Sigma}$ with (i, i') -th element $K(\|\mathbf{x}^{(i)} - \mathbf{x}^{(i')}\|)$ is a positive definite correlation matrix. Necessarily, $K(0) = 1$.

The autocovariance function of the Gaussian process Y is $\sigma^2 K_\theta$, where K_θ is the autocorrelation kernel parametrized by $\theta \in (0, +\infty)$ and defined by $K_\theta(d) = K(d/\theta)$, making $\sigma^2 \in (0, +\infty)$ the variance of $Y(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{D}$.

Fix $n \in \mathbb{Z}_+$ and a collection of n distinct points $(\mathbf{x}^{(i)})_{i \in \llbracket 1, n \rrbracket}$. Let this collection be the design set, i.e. the set of points where Y is observed. $(Y(\mathbf{x}^{(1)}), \dots, Y(\mathbf{x}^{(n)}))^\top$ is a Gaussian vector with mean vector $(f(\mathbf{x}^{(1)}), \dots, f(\mathbf{x}^{(n)}))^\top$ and covariance matrix $\sigma^2 \boldsymbol{\Sigma}_\theta$, where $\boldsymbol{\Sigma}_\theta$ denotes the $n \times n$ matrix with (i, i') -th element $K_\theta(\|\mathbf{x}^{(i)} - \mathbf{x}^{(i')}\|)$. Table 3.1 recalls Table 1.1, which provided the definition of several correlation kernels.

Kernel	$K_\theta(x)$	parameter range
Spherical ($r = 1, 2, 3$)	$\left(1 - \frac{3}{2} \left(\frac{ x }{\theta}\right) + \frac{1}{2} \left(\frac{ x }{\theta}\right)^3\right) \mathbf{1}_{\{ x \leq \theta\}}$	\emptyset
Power Exponential	$\exp\left\{-\left(\frac{ x }{\theta}\right)^q\right\}$	$q \in (0, 2]$
Rational Quadratic	$\left(1 + \left(\frac{ x }{\theta}\right)^2\right)^{-\nu}$	$\nu \in (0, +\infty)$
Matérn	$\Gamma(\nu)^{-1} 2^{1-\nu} \left(2\sqrt{\nu} \frac{ x }{\theta}\right)^\nu \mathcal{K}_\nu\left(2\sqrt{\nu} \frac{ x }{\theta}\right)$	$\nu \in (0, +\infty)$

Table 3.1 – Formulas for several correlation kernel families. The Squared Exponential kernel is the Power Exponential kernel with $q = 2$. \mathcal{K}_ν is the modified Bessel function of second kind with parameter ν [Abramowitz and Stegun, 1964](9.6.). This parametrization of the Matérn family is recommended by Handcock and Wallis [1994]. To recover the one used by Berger et al. [2001], simply replace $2\sqrt{\nu}|x|$ by $|x|$.

If p is non-zero, let \mathbf{H} denote the $n \times p$ matrix with (i, j) -th element $f_j(\mathbf{x}^{(i)})$. [Note: if $p = 0$, then we adopt the convention that any term involving \mathbf{H} can be ignored.] Then $(f(\mathbf{x}^{(1)}), \dots, f(\mathbf{x}^{(n)}))^\top = \mathbf{H}\boldsymbol{\beta}$. Denote by $\mathbf{y} = (y_1, \dots, y_n)^\top$ the observed value of the random vector $(Y(\mathbf{x}^{(1)}), \dots, Y(\mathbf{x}^{(n)}))^\top$. The likelihood function of the parameter triplet $(\boldsymbol{\beta}, \sigma^2, \theta)$ has the following expression:

$$L(\mathbf{y} | \boldsymbol{\beta}, \sigma^2, \theta) = \left(\frac{1}{2\pi\sigma^2}\right)^{\frac{n}{2}} |\boldsymbol{\Sigma}_\theta|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{H}\boldsymbol{\beta})^\top \boldsymbol{\Sigma}_\theta^{-1}(\mathbf{y} - \mathbf{H}\boldsymbol{\beta})\right\}. \quad (3.1)$$

In order for the model to be identifiable, assume that $p < n$ and that \mathbf{H} has rank p .

Berger et al. [2001] derive the reference prior corresponding to the parameter ordering $\boldsymbol{\beta} \prec (\sigma^2, \theta)$ [if $p = 0$, $\boldsymbol{\beta}$ is meaningless, so the ordering is (σ^2, θ)]. One can see [Ren et al., 2012] that the reference prior corresponding to the ordering $\boldsymbol{\beta} \prec \sigma^2 \prec \theta$ [if $p = 0$, $\sigma^2 \prec \theta$] is the same.

To express it conveniently, denote by \mathbf{Q}_θ the matrix $\mathbf{I}_n - \mathbf{H}(\mathbf{H}^\top \boldsymbol{\Sigma}_\theta^{-1} \mathbf{H})^{-1} \mathbf{H}^\top \boldsymbol{\Sigma}_\theta^{-1}$ [if $p = 0$, $\mathbf{Q}_\theta = \mathbf{I}_n$]. Also fix \mathbf{W} , an $n \times (n - p)$ matrix such that $\mathbf{W}^\top \mathbf{W} = \mathbf{I}_{n-p}$ and $\mathbf{H}^\top \mathbf{W}$ is the $p \times (n - p)$ null matrix. \mathbf{W} 's columns form an orthonormal basis of the orthogonal complement of the subspace of \mathbb{R}^n spanned by the columns of \mathbf{H} [if $p = 0$, fix \mathbf{W} as an orthogonal matrix, for instance \mathbf{I}_n].

$\boldsymbol{\beta}$ is a location parameter and $\sigma := \sqrt{\sigma^2}$ a scale parameter. Therefore, conditional on θ , the reference prior with ordering $\boldsymbol{\beta} \prec \sigma$ is proportional to $d\boldsymbol{\beta}\sigma^{-1}d\sigma$. Because the multi-parameter reference prior is invariant under parameter-by-parameter reparametrizations and $d\sigma^2/d\sigma = 2\sigma$, the reference prior with ordering $\boldsymbol{\beta} \prec \sigma^2$ is proportional to $d\boldsymbol{\beta}(\sigma^2)^{-1}d(\sigma^2)$.

Proposition 3.1. *If $p \geq 1$, after marginalizing $\boldsymbol{\beta}$ and σ^2 out, we have*

$$\begin{aligned} L(\mathbf{y}|\theta) &= \iint L(\mathbf{y}|\boldsymbol{\beta}, \sigma^2, \theta)/\sigma^2 d\boldsymbol{\beta} d\sigma^2 \\ &= \left(\frac{2\pi^{\frac{n-p}{2}}}{\Gamma(\frac{n-p}{2})}\right)^{-1} |\boldsymbol{\Sigma}_\theta^{-1}|^{\frac{1}{2}} \left|\mathbf{H}^\top \boldsymbol{\Sigma}_\theta^{-1} \mathbf{H}\right|^{-\frac{1}{2}} (\mathbf{y}^\top \boldsymbol{\Sigma}_\theta^{-1} \mathbf{Q}_\theta \mathbf{y})^{\frac{n-p}{2}}. \end{aligned} \quad (3.2)$$

Alternatively, the integrated likelihood with $p \geq 1$ can also be written

$$L(\mathbf{y}|\theta) = \left(\frac{2\pi^{\frac{n-p}{2}}}{\Gamma(\frac{n-p}{2})}\right)^{-1} \left|\mathbf{H}^\top \mathbf{H}\right|^{\frac{1}{2}} \left|\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W}\right|^{-\frac{1}{2}} \left(\mathbf{y}^\top \mathbf{W} (\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{y}\right)^{\frac{n-p}{2}}. \quad (3.3)$$

If $p = 0$, the integrated likelihood is simply

$$L(\mathbf{y}|\theta) = \int L(\mathbf{y}|\sigma^2, \theta)/\sigma^2 d\sigma^2 = \left(\frac{2\pi^{\frac{n}{2}}}{\Gamma(\frac{n}{2})} \right)^{-1} |\boldsymbol{\Sigma}_\theta^{-1}|^{\frac{1}{2}} |(\mathbf{y}^\top \boldsymbol{\Sigma}_\theta^{-1} \mathbf{y})^{\frac{n}{2}}|. \quad (3.4)$$

Proof. The result for $p = 0$ and the first result for $p \geq 1$ are from Berger et al. [2001]. Lemma 3.10 yields that

$$\mathbf{W} \left(\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} \right)^{-1} \mathbf{W}^\top = \boldsymbol{\Sigma}_\theta^{-1} \mathbf{Q}_\theta. \quad (3.5)$$

So all that remains to be proved is that $|\boldsymbol{\Sigma}_\theta| = \left| \mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} \right| \left| \mathbf{H}^\top \mathbf{H} \right|^{-1} \left| \mathbf{H}^\top \boldsymbol{\Sigma}_\theta^{-1} \mathbf{H} \right|^{-1}$. Choose an $n \times p$ matrix \mathbf{P} with columns forming an orthonormal basis of the subspace of \mathbb{R}^n spanned by the columns of \mathbf{H} . $(\mathbf{W}\mathbf{P})$ is therefore an $n \times n$ orthogonal matrix, so $|\boldsymbol{\Sigma}_\theta| = |(\mathbf{W}\mathbf{P})^\top \boldsymbol{\Sigma}_\theta (\mathbf{W}\mathbf{P})|$. Using Schur's complement, we have

$$|\boldsymbol{\Sigma}_\theta| = \left| \mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} \right| \left| \mathbf{P}^\top \boldsymbol{\Sigma}_\theta \left(\mathbf{I}_n - \mathbf{W} \left(\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} \right)^{-1} \mathbf{W}^\top \boldsymbol{\Sigma}_\theta \right) \mathbf{P} \right|. \quad (3.6)$$

Lemma 3.10 again yields the result. \square

In the following proposition, the first assertion is from Ren et al. [2012].

Proposition 3.2. *The reference prior with ordering $\boldsymbol{\beta} \prec \sigma^2 \prec \theta$ is $\pi(\boldsymbol{\beta}, \sigma^2, \theta) \propto (\sigma^2)^{-1} \pi(\theta)$, where*

$$\pi(\theta) \propto \sqrt{\text{Tr} \left[\left\{ \left(\frac{d}{d\theta} \boldsymbol{\Sigma}_\theta \right) \boldsymbol{\Sigma}_\theta^{-1} \mathbf{Q}_\theta \right\}^2 \right]} - \frac{1}{n-p} \left[\text{Tr} \left\{ \left(\frac{d}{d\theta} \boldsymbol{\Sigma}_\theta \right) \boldsymbol{\Sigma}_\theta^{-1} \mathbf{Q}_\theta \right\} \right]^2. \quad (3.7)$$

Denoting $\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W}$ by $\boldsymbol{\Sigma}_\theta^{\mathbf{W}}$, $\pi(\theta)$ can also be written as:

$$\pi(\theta) \propto \sqrt{\text{Tr} \left[\left\{ \left(\frac{d}{d\theta} \boldsymbol{\Sigma}_\theta^{\mathbf{W}} \right) \left(\boldsymbol{\Sigma}_\theta^{\mathbf{W}} \right)^{-1} \right\}^2 \right]} - \frac{1}{n-p} \left[\text{Tr} \left\{ \left(\frac{d}{d\theta} \boldsymbol{\Sigma}_\theta^{\mathbf{W}} \right) \left(\boldsymbol{\Sigma}_\theta^{\mathbf{W}} \right)^{-1} \right\} \right]^2. \quad (3.8)$$

Proof. To obtain the second assertion, all we need to do is recognize Example 3 from Chapter 2. This can be done after noticing that $\mathbf{W}^\top \mathbf{y}$ plays the role of \mathbf{y} and $\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W}$ the role of $\boldsymbol{\Sigma}_\theta$. Because $\mathbf{W}^\top \mathbf{y} \in \mathbb{R}^{n-p}$, $n-p$ plays the role of n .

The first assertion is then a consequence of Lemma 3.10. \square

3.3 Smoothness of the correlation kernel

Lemma 2 of Berger et al. [2001] requires that correlation kernel and design set should be such that $\boldsymbol{\Sigma}_\theta = \mathbf{1}\mathbf{1}^\top + g_0(\theta)\mathbf{D} + \mathbf{R}_0(\theta)$, where $\mathbf{1}$ is the vector with n entries all equal to 1, $g_0(\theta)$ is a real-valued function such that $\lim_{\theta \rightarrow +\infty} g_0(\theta) = 0$, \mathbf{D} is a fixed nonsingular matrix and \mathbf{R}_0 is a mapping from $(0, +\infty)$ to the set of $n \times n$ real matrices \mathcal{M}_n such that $\lim_{\theta \rightarrow +\infty} \left\| \frac{1}{g_0(\theta)} \mathbf{R}_0(\theta) \right\| = 0$.

What makes this assumption restrictive is the condition that \mathbf{D} should be nonsingular, because it holds for rough correlation kernels only. For instance, as was noted by Paulo [2005], it does not hold for the Squared Exponential correlation kernel.

For a given correlation kernel K , \mathbf{D} is typically a matrix proportional to the matrix with entries $\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^q$, where q depends on the smoothness of the correlation kernel but should in any case belong to the interval $(0, 2]$. This is because $K(s) - K(0)$ is equivalent to a constant times s^q when $s \rightarrow 0+$.

Schoenberg [1937] gives the following result (Theorem 4 in the original paper):

Theorem 3.3. *If $q \in (0, 2)$, the quadratic form $\boldsymbol{\xi} \in \mathbb{R}^n \mapsto \sum_{i,j=0}^n \|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^q \xi_i \xi_j$ is nonsingular and its canonical representation contains one positive and n negative squares.*

This means that if the correlation kernel is rough enough to have $q \in (0, 2)$, the assumption that \mathbf{D} is nonsingular is reasonable.

Corollary 3.4. *The $n \times n$ matrix with entries $\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^q$ with $q \in (0, 2)$ is nonsingular and has one positive eigenvalue and n negative eigenvalues.*

The picture is dramatically different when the correlation kernel K is smooth enough to have $q = 2$. This happens as soon as K is twice continuously differentiable. Gower [1985]’s Theorem 6 implies the following results:

Theorem 3.5. *If d is the dimension of E_d , the smallest Euclidean subspace containing all points in the design set, then the $n \times n$ matrix with entries $\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2$ has rank:*

- (a) $d + 1$ (one positive eigenvalue, d negative eigenvalues, any other eigenvalue null) if all points in the design set lie on the surface of a hypersphere of E_d ;
- (b) $d + 2$ (one positive eigenvalue, $d + 1$ negative eigenvalues, any other eigenvalue null) otherwise.

Corollary 3.6. *The $n \times n$ matrix with entries $\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2$ has rank lower or equal to $r + 2$.*

For all practical purposes, n is much greater than r , so the matrix \mathbf{D} is singular when $q = 2$.

Let us review the values of q for correlation kernels listed in Table 3.1. Matérn correlation kernels [Matérn, 1986] [Handcock and Stein, 1993] with smoothness parameter ν have $q = 2 \min(1, \nu)$, thus for $0 < \nu < 1$, $0 < q < 2$ but for $\nu \geq 1$, $q = 2$. Spherical correlation kernels [Wackernagel, 1995] have $q = 1$. Power Exponential kernels [De Oliveira et al., 1997] have q equal to their power. This means that all Power Exponential kernels except the Squared Exponential correlation kernel have $0 < q < 2$. In particular, the Exponential kernel (which is also the Matérn kernel with smoothness $\nu = 1/2$) has $q = 1$, but the Squared Exponential kernel has $q = 2$. Rational Quadratic kernels [Yaglom, 1987] have $q = 2$. For easy reference, the review is summarized in Table 3.2.

This review justifies the claim in the abstract that the Squared Exponential kernel, Matérn kernels with smoothness $\nu \geq 1$ and Rational Quadratic kernels require a proof of the reference posterior’s propriety.

Kernel	$g_0(\theta)$	$\ \mathbf{R}_0(\theta)\ $	q	\mathbf{D} nonsingular*
Spherical ($r = 1, 2, 3$)	$-3/2\theta^{-1}$	$O(\theta^{-3})$	1	yes
Power Expon. ($q < 2$)	$-\theta^{-q}$	$O(\theta^{-2q})$	q	yes
Squared Exponential	$-\theta^{-2}$	$O(\theta^{-4})$	2	no
Rational Quadratic	$-\nu\theta^{-2}$	$O(\theta^{-4})$	2	no
Matérn ($\nu < 1$)	$\Gamma(-\nu)\nu^\nu\Gamma(\nu)^{-1}\theta^{-2\nu}$	$O(\theta^{-2})$	2ν	yes
Matérn ($\nu = 1$)	$-2\theta^{-2}\log(\theta)$	$O(\theta^{-2})$	2	no
Matérn ($1 < \nu < 2$)	$-\Gamma(\nu - 1)\nu\Gamma(\nu)^{-1}\theta^2$	$O(\theta^{-2\nu})$	2	no
Matérn ($\nu = 2$)	$-2\theta^2$	$O(\theta^{-4\log(\theta)})$	2	no
Matérn ($\nu > 2$)	$-\Gamma(\nu - 1)\nu\Gamma(\nu)^{-1}\theta^2$	$O(\theta^{-4})$	2	no

Table 3.2 – Summary of the results of Section 3.3. *Answer given assuming $n > r + 2$.

3.4 Propriety of the reference posterior distribution

Berger et al. [2001] show that the reference posterior distribution on β and σ^2 conditionally to θ is proper.

In this section, we prove that the joint reference posterior distribution is proper for Matérn kernels with smoothness $\nu \geq 1$, Rational Quadratic kernels and the Squared Exponential kernel.

Proposition 3.7. *For Matérn kernels with smoothness $\nu \geq 1$, for Rational Quadratic kernels with parameter $\nu > 0$ and for the Squared Exponential kernel, the “marginal” reference prior distribution $\pi(\theta)$ defined by Proposition 3.2 has the following behavior.*

1. When $\theta \rightarrow 0$,

$$\pi(\theta) = \begin{cases} o(1) & \text{for Matérn kernels and the Squared Exponential kernel;} \\ O(\theta^{2\nu-1}) & \text{for Rational Quadratic kernels.} \end{cases} \quad (3.9)$$

2. When $\theta \rightarrow +\infty$,

$$\pi(\theta) = \begin{cases} O(\theta^{-1}) & \text{for Matérn kernels;} \\ o(1) & \text{for Rational Quadratic kernels;} \\ O(\theta) & \text{for the Squared Exponential kernel.} \end{cases} \quad (3.10)$$

Proof. Denoting any of these kernels by K , K is continuously differentiable.

If K is Squared Exponential, $\lim_{\theta \rightarrow 0} \frac{d}{d\theta} K(1/\theta) = 0$. This also holds if K is Matérn with smoothness $\nu \geq 1$ (see Abramowitz and Stegun [1964] 9.6.28. and 9.7.2.). If K is Rational Quadratic with parameter $\nu > 0$, $\frac{d}{d\theta} K(1/\theta) \underset{\theta \rightarrow 0}{\sim} 2\nu\theta^{2\nu-1}$. Moreover, Σ_θ converges to \mathbf{I}_n when $\theta \rightarrow 0$, so its inverse does too. The first assertion follows from these facts.

The second assertion is proved by combining Lemma 3.12 with Lemma 3.20/3.21/3.22 for Matérn, Rational Quadratic and Squared Exponential kernels respectively. \square

Let $v_1(\theta) \geq \dots \geq v_{n-p}(\theta) > 0$ be the ordered eigenvalues of $\mathbf{W}^\top \Sigma_\theta \mathbf{W}$.

Lemma 3.8. *For Rational Quadratic and Squared Exponential kernels and for Matérn kernels with smoothness $\nu \geq 1$, there exists a hyperplane \mathcal{H} of \mathbb{R}^n such that for every $\mathbf{y} \in \mathbb{R}^n \setminus \mathcal{H}$, when $\theta \rightarrow +\infty$:*

$$\left(\mathbf{y}^\top \mathbf{W} \left(\mathbf{W}^\top \Sigma_\theta \mathbf{W} \right)^{-1} \mathbf{W}^\top \mathbf{y} \right)^{-1} = O(v_{n-p}(\theta)). \quad (3.11)$$

The proof of this lemma can be found in Appendix 3.D. Combined with Equation (3.3), it implies that if the observation vector \mathbf{y} belongs to $\mathbb{R}^n \setminus \mathcal{H}$, then

$$L(\mathbf{y}|\theta)^2 = \prod_{i=1}^{n-p} \frac{O(v_{n-p}(\theta))}{v_i(\theta)} = O(1) \quad \text{when } \theta \rightarrow +\infty. \quad (3.12)$$

In the following, when \mathbf{y} belongs to $\mathbb{R}^n \setminus \mathcal{H}$, we write that “ \mathbf{y} looks nondegenerate”. This terminology relies on the intuition that if the observation were to take some values within \mathcal{H} , it would be better explained by a degenerate Gaussian model. The most compelling example is that of a constant observation vector, for which the Kriging model would be grossly inappropriate.

Theorem 3.9. *For Matérn kernels with noninteger smoothness $\nu > 1$, for Rational Quadratic kernels and for the Squared Exponential kernel, regardless of the design set and of the mean function space, if \mathbf{y} looks nondegenerate, then the reference posterior distribution $\pi(\theta|\mathbf{y})$ is proper.*

Proof. The first assertion of Proposition 3.7 implies the reference prior $\pi(\theta)$ is integrable in the neighborhood of 0. Furthermore, when $\theta \rightarrow 0$, $\Sigma_\theta \rightarrow \mathbf{I}_n$ so the reference posterior $\pi(\theta|\mathbf{y})$ is integrable in the neighborhood of 0 as well.

All that remains to be proved is therefore that the reference posterior is integrable in the neighborhood of $+\infty$. In the following $\theta \rightarrow +\infty$, so we rely on the asymptotic expansion of Σ_θ , which is detailed in Appendix 3.D.

The proof is somewhat trickier for Matérn kernels with integer smoothness, so we tackle this case at the end. Until further notice, assume the kernel is Rational Quadratic, Squared Exponential or Matérn with noninteger smoothness $\nu > 1$.

For Rational Quadratic and Squared Exponential (resp. Matérn with noninteger smoothness parameter $\nu > 1$) kernels, Appendix 3.D (resp. Appendix 3.D) shows how $\mathbf{W}^\top \Sigma_\theta \mathbf{W}$ can be decomposed as

$$\mathbf{W}^\top \Sigma_\theta \mathbf{W} = g(\theta) \left(\mathbf{W}^\top \mathbf{D} \mathbf{W} + g^*(\theta) \mathbf{W}^\top \mathbf{D}^* \mathbf{W} + \mathbf{R}_g(\theta) \right), \quad (3.13)$$

where

- g is a differentiable function;
- $g^*(\theta) = \theta^{-2l}$ with $l \in (0, +\infty)$ (actually, if the kernel is Rational Quadratic or Squared Exponential, $l \in \mathbb{Z}_+$);
- \mathbf{R}_g is a differentiable mapping from $(0, +\infty)$ to \mathcal{M}_n such that $\|\mathbf{R}_g(\theta)\| = o(\theta^{-2l})$;
- \mathbf{D} and \mathbf{D}^* are both fixed symmetric matrices;
- $\mathbf{W}^\top \mathbf{D} \mathbf{W}$ is non-null;
- either $\mathbf{W}^\top \mathbf{D}^* \mathbf{W}$ is non-null or $\mathbf{W}^\top \mathbf{D} \mathbf{W}$ is nonsingular.

Lemma 3.15 implies that one of the following is true:

1. When θ is large enough, $\mathbf{W}^\top \mathbf{D} \mathbf{W} + g^*(\theta) \mathbf{W}^\top \mathbf{D}^* \mathbf{W}$ is nonsingular. This case can be further decomposed in the following subcases:
 - a) $\mathbf{W}^\top \mathbf{D} \mathbf{W}$ is nonsingular;
 - b) $\mathbf{W}^\top \mathbf{D} \mathbf{W}$ is singular, but $\mathbf{W}^\top \mathbf{D} \mathbf{W} + g^*(\theta) \mathbf{W}^\top \mathbf{D}^* \mathbf{W}$ is nonsingular when θ is large enough.

2. The vector space $\text{Ker}(\mathbf{W}^\top \mathbf{D}^* \mathbf{W}) \cap \text{Ker}(\mathbf{W}^\top \mathbf{D} \mathbf{W})$ is non-trivial.

Let us differentiate $\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W}$:

$$\frac{d}{d\theta} \mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} = \frac{g'(\theta)}{g(\theta)} \mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} + g(\theta) \left(g^{*\prime}(\theta) \mathbf{W}^\top \mathbf{D}^* \mathbf{W} + \frac{d}{d\theta} \mathbf{R}_g(\theta) \right). \quad (3.14)$$

We can show that $\|\frac{d}{d\theta} \mathbf{R}_g(\theta)\| = o(g^{*\prime}(\theta))$. This is due to Equation (3.72) for Rational Quadratic and Squared Exponential kernels, and to Equation (3.75) for Matérn kernels with noninteger smoothness.

Lemma 3.11 shows that $\frac{d}{d\theta} \mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W}$ can be replaced by $g(\theta) \left(g^{*\prime}(\theta) \mathbf{W}^\top \mathbf{D}^* \mathbf{W} + \frac{d}{d\theta} \mathbf{R}_g(\theta) \right)$ in Equation (3.8): $\pi(\theta) \propto w(\theta)$, where

$$w(\theta)^2 := \text{Tr} \left[\left\{ g(\theta) \left(g^{*\prime}(\theta) \mathbf{W}^\top \mathbf{D}^* \mathbf{W} + \frac{d}{d\theta} \mathbf{R}_g(\theta) \right) \left(\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} \right)^{-1} \right\}^2 \right] - \frac{1}{n-p} \left[\text{Tr} \left\{ g(\theta) \left(g^{*\prime}(\theta) \mathbf{W}^\top \mathbf{D}^* \mathbf{W} + \frac{d}{d\theta} \mathbf{R}_g(\theta) \right) \left(\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} \right)^{-1} \right\} \right]^2. \quad (3.15)$$

We have $w(\theta) \leq \tilde{w}(\theta)$, where

$$\tilde{w}(\theta) := \sqrt{\text{Tr} \left[\left\{ g(\theta) \left(g^{*\prime}(\theta) \mathbf{W}^\top \mathbf{D}^* \mathbf{W} + \frac{d}{d\theta} \mathbf{R}_g(\theta) \right) \left(\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} \right)^{-1} \right\}^2 \right]}. \quad (3.16)$$

A specific asymptotic analysis is required in each case. This study is conducted in Appendix 3.E. We summarize the results in Table 3.3.

Case	Kernels	$\pi(\theta)$	$L(\mathbf{y} \theta)$
1.(a)	Matérn ($\nu \in [1, +\infty) \setminus \mathbb{Z}_+$), RQ, SE	$O(\theta^{-2l-1})$	$O(1)$
1.(b)	Matérn ($\nu \in [1, +\infty) \setminus \mathbb{Z}_+$), RQ, SE	$O(\theta^{-1})$	$O(\theta^{-l})$
2.	Matérn ($\nu \in [1, +\infty) \setminus \mathbb{Z}_+$)	$O(\theta^{-1})$	$O(\theta^{-l})$
2.	RQ, SE (usual case)	$O(\theta)$	$O(\theta^{-3})$
2.	RQ, SE (special case)	$O(\theta^{-1})$	$O(\theta^{-1})$

Table 3.3 – Asymptotic upper bounds for reference prior $\pi(\theta)$ and likelihood $L(\mathbf{y}|\theta)$ for Rational Quadratic (RQ) and Squared Exponential (SE) kernels and Matérn kernels with noninteger smoothness $\nu > 1$ in all three cases. The proof in Appendix 3.E shows that for Rational Quadratic and Squared Exponential kernels, case 2. can be split in two subcases (“usual” and “special”).

The posterior distribution resulting from the reference prior is proper in all cases.

Matérn kernels with integer smoothness are dealt with in Appendix 3.E. □

3.5 Conclusion

The main result of this chapter is Theorem 3.9, which ensures that the reference prior leads to a proper posterior distribution for a large class of smooth kernels. This class contains the Squared Exponential correlation kernel as well as the important Matérn family [Stein, 1999] with smoothness parameter $\nu \geq 1$. Rational Quadratic kernels, whose usage is less widespread are also included within this class.

Berger et al. [2001] proved this result for a class of rough correlation kernels. This class includes the complementary set of the Matérn family – kernels with smoothness parameter $\nu < 1$ – as well as all other Power Exponential kernels. Spherical kernels, which are mostly used in the field of geostatistics also belong to this class.

Combining Theorem 3.9 with the results from Berger et al. [2001], one can appreciate how polyvalent the reference prior is, insofar as it is able to adapt to very different correlation kernels and always leads to a proper posterior. No *ad-hoc* technique is required to derive useable inference, so this approach seems to be flawless from a Bayesian point of view when no explicit prior information is available. Even when explicit prior information is available, following Druilhet and Marin [2007], it can be used to derive Maximum A Posteriori (MAP) estimates or High Probability Density (HPD) sets that are invariant under reparametrization.

Appendix 3.A Algebraic facts

Lemma 3.10. *Let a and b be positive integers and let Σ be a nonsingular symmetric $(a + b) \times (a + b)$ matrix. Then, for any $(a + b) \times a$ matrix \mathbf{A} with rank a and any $(a + b) \times b$ matrix \mathbf{B} with rank b such that $\mathbf{A}^\top \mathbf{B}$ is the null $a \times b$ matrix,*

$$\mathbf{B} \left(\mathbf{B}^\top \Sigma \mathbf{B} \right)^{-1} \mathbf{B}^\top = \Sigma^{-1} \left(\mathbf{I}_{a+b} - \mathbf{A} \left(\mathbf{A}^\top \Sigma^{-1} \mathbf{A} \right)^{-1} \mathbf{A}^\top \Sigma^{-1} \right). \quad (3.17)$$

Proof. Notice that both matrices have the same kernel, namely the subspace of \mathbb{R}^{a+b} spanned by \mathbf{A} . Indeed, because \mathbf{B} has full column rank and $\left(\mathbf{B}^\top \Sigma \mathbf{B} \right)^{-1}$ is nonsingular, the left matrix has the same kernel as \mathbf{B}^\top . Besides, the a -dimensional subspace of \mathbb{R}^{a+b} spanned by \mathbf{A} is included in this kernel. So because the rank of \mathbf{B}^\top is b , its kernel has dimension a and the inclusion is an equality.

Similarly, because Σ^{-1} is nonsingular, the right matrix has the same kernel as $\mathbf{I}_{a+b} - \mathbf{A} \left(\mathbf{A}^\top \Sigma^{-1} \mathbf{A} \right)^{-1} \mathbf{A}^\top \Sigma^{-1}$. Moreover, because the image of $\mathbf{A} \left(\mathbf{A}^\top \Sigma^{-1} \mathbf{A} \right)^{-1} \mathbf{A}^\top \Sigma^{-1}$ is included within the image of \mathbf{A} , its dimension is lower or equal to a . The image of \mathbf{I}_{a+b} on the other hand has dimension $a + b$, so the image of $\mathbf{I}_{a+b} - \mathbf{A} \left(\mathbf{A}^\top \Sigma^{-1} \mathbf{A} \right)^{-1} \mathbf{A}^\top \Sigma^{-1}$ has dimension greater or equal to b and therefore its kernel has dimension lower or equal to a . Now, a simple computation shows that the a -dimensional subspace of \mathbb{R}^{a+b} spanned by \mathbf{A} is included in the kernel, so it is in fact equal to the kernel.

Besides, for any $\mathbf{z} \in \mathbb{R}^b$,

$$\mathbf{B} \left(\mathbf{B}^\top \Sigma \mathbf{B} \right)^{-1} \mathbf{B}^\top (\Sigma \mathbf{B} \mathbf{z}) = \mathbf{B} \mathbf{z}; \quad (3.18)$$

$$\Sigma^{-1} \left(\mathbf{I}_{a+b} - \mathbf{A} \left(\mathbf{A}^\top \Sigma^{-1} \mathbf{A} \right)^{-1} \mathbf{A}^\top \Sigma^{-1} \right) (\Sigma \mathbf{B} \mathbf{z}) = \mathbf{B} \mathbf{z}. \quad (3.19)$$

So both matrices act the same way on the subspace spanned by $\Sigma \mathbf{B}$, which is supplementary to their common kernel, hence the equality. \square

Lemma 3.11. *Let m be a positive integer, Σ be a nonsingular $m \times m$ matrix, and \mathbf{A} and \mathbf{B} be $m \times m$ matrices. If there exists a real number t such that*

$$\mathbf{A} = t\Sigma + \mathbf{B}, \quad (3.20)$$

then

$$\text{Tr} \left[\left\{ \mathbf{A} \Sigma^{-1} \right\}^2 \right] - \frac{1}{m} \left[\text{Tr} \left\{ \mathbf{A} \Sigma^{-1} \right\} \right]^2 = \text{Tr} \left[\left\{ \mathbf{B} \Sigma^{-1} \right\}^2 \right] - \frac{1}{m} \left[\text{Tr} \left\{ \mathbf{B} \Sigma^{-1} \right\} \right]^2. \quad (3.21)$$

Proof. The lemma follows from a direct calculation:

$$\text{Tr} \left[\mathbf{A} \Sigma^{-1} \right] = \text{Tr} \left[\mathbf{B} \Sigma^{-1} \right] + tm \quad (3.22)$$

$$\text{Tr} \left[\left\{ \mathbf{A} \Sigma^{-1} \right\}^2 \right] = \text{Tr} \left[\left\{ \mathbf{B} \Sigma^{-1} \right\}^2 \right] + 2t \text{Tr} \left[\mathbf{B} \Sigma^{-1} \right] + t^2 m \quad (3.23)$$

\square

Lemma 3.12. *Let $m > a$ be positive integers, Σ be an $m \times m$ symmetric positive definite matrix, Σ' be an $m \times m$ symmetric matrix and \mathbf{A} be an $m \times a$ matrix with rank a . Denote by \mathbf{Q} the matrix $\mathbf{I}_m - \mathbf{A} \left(\mathbf{A}^\top \Sigma^{-1} \mathbf{A} \right)^{-1} \mathbf{A}^\top \Sigma^{-1}$. Then, if there exist $t_1 \in \mathbb{R}$ and $t_2 \in [0, +\infty)$ such that the matrix $\mathbf{F} := t_1 \Sigma - \Sigma'$ is positive semi-definite and verifies $\forall \xi \in \mathbb{R}^m \quad \xi^\top \mathbf{F} \xi \leq t_2 \xi^\top \Sigma \xi$, then*

$$\sqrt{\text{Tr} [(\Sigma' \Sigma^{-1} \mathbf{Q})^2] - \frac{1}{m-a} \text{Tr} [\Sigma' \Sigma^{-1} \mathbf{Q}]^2} \leq (m-a)t_2 \quad (3.24)$$

Proof. Let \mathbf{B} be an $m \times (m-a)$ matrix with rank $m-a$ such that $\mathbf{A}^\top \mathbf{B}$ is the null $a \times (m-a)$ matrix. Such a matrix \mathbf{B} can for instance be constructed by computing a Singular Value Decomposition (SVD) of \mathbf{A} : $\mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^\top$. In this decomposition, \mathbf{U} and \mathbf{V} are orthogonal matrices of size $m \times m$ and $a \times a$ respectively, and \mathbf{S} is an $m \times a$ matrix whose only non-null entries are on the main diagonal. Therefore the last $m-a$ rows of \mathbf{S} are filled with zeros. So define \mathbf{B} as the $m \times (m-a)$ matrix formed by the last $m-a$ columns of \mathbf{U} .

By applying Lemma 3.10, we obtain that $\Sigma^{-1} \mathbf{Q} = \mathbf{B} \left(\mathbf{B}^\top \Sigma \mathbf{B} \right)^{-1} \mathbf{B}^\top$.

Because of the properties of the trace, this implies

$$\text{Tr} [\Sigma' \Sigma^{-1} \mathbf{Q}] = \text{Tr} \left[\mathbf{B}^\top \Sigma' \mathbf{B} \left(\mathbf{B}^\top \Sigma \mathbf{B} \right)^{-1} \right] \quad (3.25)$$

$$\text{Tr} [(\Sigma' \Sigma^{-1} \mathbf{Q})^2] = \text{Tr} \left[\left\{ \mathbf{B}^\top \Sigma' \mathbf{B} \left(\mathbf{B}^\top \Sigma \mathbf{B} \right)^{-1} \right\}^2 \right]. \quad (3.26)$$

Similarly, we have

$$\text{Tr} [\mathbf{F} \Sigma^{-1} \mathbf{Q}] = \text{Tr} \left[\mathbf{B}^\top \mathbf{F} \mathbf{B} \left(\mathbf{B}^\top \Sigma \mathbf{B} \right)^{-1} \right] \quad (3.27)$$

$$\text{Tr} [(\mathbf{F} \Sigma^{-1} \mathbf{Q})^2] = \text{Tr} \left[\left\{ \mathbf{B}^\top \mathbf{F} \mathbf{B} \left(\mathbf{B}^\top \Sigma \mathbf{B} \right)^{-1} \right\}^2 \right]. \quad (3.28)$$

Because $\mathbf{B}^\top \mathbf{F} \mathbf{B} = t_1 \mathbf{B}^\top \Sigma \mathbf{B} - \mathbf{B}^\top \Sigma' \mathbf{B}$, Lemma 3.11 implies

$$\begin{aligned} & \text{Tr} \left[\left\{ \mathbf{B}^\top \Sigma' \mathbf{B} \left(\mathbf{B}^\top \Sigma \mathbf{B} \right)^{-1} \right\}^2 \right] - \frac{1}{m-a} \text{Tr} \left[\mathbf{B}^\top \Sigma' \mathbf{B} \left(\mathbf{B}^\top \Sigma \mathbf{B} \right)^{-1} \right]^2 \\ &= \text{Tr} \left[\left\{ \mathbf{B}^\top \mathbf{F} \mathbf{B} \left(\mathbf{B}^\top \Sigma \mathbf{B} \right)^{-1} \right\}^2 \right] - \frac{1}{m-a} \text{Tr} \left[\mathbf{B}^\top \mathbf{F} \mathbf{B} \left(\mathbf{B}^\top \Sigma \mathbf{B} \right)^{-1} \right]^2. \end{aligned} \quad (3.29)$$

Combining the 5 equations above yields

$$\text{Tr} [(\Sigma' \Sigma^{-1} \mathbf{Q})^2] - \frac{1}{m-a} \text{Tr} [\Sigma' \Sigma^{-1} \mathbf{Q}]^2 = \text{Tr} [(\mathbf{F} \Sigma^{-1} \mathbf{Q})^2] - \frac{1}{m-a} \text{Tr} [\mathbf{F} \Sigma^{-1} \mathbf{Q}]^2. \quad (3.30)$$

An elementary computation shows that $\Sigma^{-1} \mathbf{Q} = \mathbf{Q}^\top \Sigma^{-1} \mathbf{Q}$.

Consider the Cholesky decomposition $\Sigma =: \mathbf{L} \mathbf{L}^\top$.

Then $\Sigma^{-1} \mathbf{Q} = \mathbf{Q}^\top \Sigma^{-1} \mathbf{Q} = \mathbf{Q}^\top (\mathbf{L}^{-1})^\top \mathbf{L}^{-1} \mathbf{Q}$.

$$\begin{aligned} \text{Tr} \left[(\mathbf{F}\boldsymbol{\Sigma}^{-1}\mathbf{Q})^2 \right] &= \text{Tr} \left[\left(\mathbf{F}\mathbf{Q}^\top (\mathbf{L}^{-1})^\top \mathbf{L}^{-1}\mathbf{Q} \right)^2 \right] = \text{Tr} \left[\left(\mathbf{L}^{-1}\mathbf{Q}\mathbf{F}\mathbf{Q}^\top (\mathbf{L}^{-1})^\top \right)^2 \right] \\ &\leq \text{Tr} \left[\mathbf{L}^{-1}\mathbf{Q}\mathbf{F}\mathbf{Q}^\top (\mathbf{L}^{-1})^\top \right]^2 = \text{Tr} \left[\mathbf{F}\boldsymbol{\Sigma}^{-1}\mathbf{Q} \right]^2. \end{aligned} \quad (3.31)$$

The inequality holds because $\mathbf{L}^{-1}\mathbf{Q}\mathbf{F}\mathbf{Q}^\top (\mathbf{L}^{-1})^\top$ is a symmetric positive semi-definite matrix.

Let $(\boldsymbol{\xi}^i)_{1 \leq i \leq m}$ be a basis of unit eigenvectors of $\boldsymbol{\Sigma}^{-1}\mathbf{Q}$ such that for every integer $i \in \llbracket 1, m \rrbracket \setminus \llbracket 1, m-a \rrbracket$, $\boldsymbol{\xi}^i$ belongs to the kernel of $\boldsymbol{\Sigma}^{-1}\mathbf{Q}$. Indeed, because $\boldsymbol{\Sigma}^{-1}\mathbf{Q} = \mathbf{B} \left(\mathbf{B}^\top \boldsymbol{\Sigma} \mathbf{B} \right)^{-1} \mathbf{B}^\top$, this kernel has the same dimension as the kernel of \mathbf{B}^\top : a .

Denoting by $(s_i)_{1 \leq i \leq m}$ the family of the eigenvalues corresponding to the family of eigenvectors $(\boldsymbol{\xi}^i)_{1 \leq i \leq m}$, we have for every integer $i \in \llbracket 1, m-a \rrbracket$ $s_i \neq 0$ and

$$\begin{aligned} (\boldsymbol{\xi}^i)^\top \boldsymbol{\Sigma} \boldsymbol{\xi}^i &= s_i^{-2} \left\{ (\boldsymbol{\xi}^i)^\top \mathbf{Q}^\top \boldsymbol{\Sigma}^{-1} \right\} \boldsymbol{\Sigma} \left\{ \boldsymbol{\Sigma}^{-1} \mathbf{Q} \boldsymbol{\xi}^i \right\} \\ &= s_i^{-2} (\boldsymbol{\xi}^i)^\top \mathbf{Q}^\top \boldsymbol{\Sigma}^{-1} \mathbf{Q} \boldsymbol{\xi}^i \\ &= s_i^{-2} (\boldsymbol{\xi}^i)^\top \boldsymbol{\Sigma}^{-1} \mathbf{Q} \boldsymbol{\xi}^i \\ &= s_i^{-1}. \end{aligned} \quad (3.32)$$

This implies the third equality below:

$$\text{Tr} \left[\mathbf{F}\boldsymbol{\Sigma}^{-1}\mathbf{Q} \right] = \sum_{i=1}^m (\boldsymbol{\xi}^i)^\top \mathbf{F}\boldsymbol{\Sigma}^{-1}\mathbf{Q}\boldsymbol{\xi}^i = \sum_{i=1}^{m-a} s_i (\boldsymbol{\xi}^i)^\top \mathbf{F}\boldsymbol{\xi}^i = \sum_{i=1}^{m-a} \frac{(\boldsymbol{\xi}^i)^\top \mathbf{F}\boldsymbol{\xi}^i}{(\boldsymbol{\xi}^i)^\top \boldsymbol{\Sigma} \boldsymbol{\xi}^i} \leq (m-a)t_2. \quad (3.33)$$

Equations (3.30) and (3.31) yield the result. □

Entire series

Lemma 3.13. *Let $(\mathbf{D}_k)_{k \in \mathbb{N}}$ be a sequence of matrices of the same size. If $\sum_{k \in \mathbb{N}} \mathbf{D}_k$ exists and its kernel is the trivial vector space, then there exists a nonnegative integer N such that $\bigcap_{k=0}^N \text{Ker } \mathbf{D}_k$ is the trivial vector space.*

Proof. Assume the sum $\sum_{k \in \mathbb{N}} \mathbf{D}_k$ exists and its kernel is the trivial vector space. Consider the sequence $(d(n))_{n \in \mathbb{N}}$ where for every nonnegative integer n , $d(n)$ is the dimension of $\bigcap_{k=0}^n \text{Ker } \mathbf{D}^{(k)}$. $(d(n))_{n \in \mathbb{N}}$ is a nonincreasing sequence of nonnegative integers, so it is convergent. If its limit is strictly greater than 0, then for every nonnegative integer n , there exists a unit vector \mathbf{v}_n that belongs to $\bigcap_{k=0}^n \text{Ker } \mathbf{D}^{(k)}$. Because the unit sphere is compact, there exists an increasing mapping $\phi : \mathbb{N} \rightarrow \mathbb{N}$ such that the subsequence $(\mathbf{v}_{\phi(n)})_{n \in \mathbb{N}}$ converges to a limit \mathbf{v} such that $\|\mathbf{v}\| = 1$. Besides, for every pair of nonnegative integers $n \leq n'$, $\mathbf{v}_{\phi(n')} \in \bigcap_{k=0}^{\phi(n)} \text{Ker } \mathbf{D}^{(k)}$. Given this set is closed, the limit \mathbf{v} also belongs to $\bigcap_{k=0}^{\phi(n)} \text{Ker } \mathbf{D}^{(k)}$. So for every nonnegative integer k , $\mathbf{v} \in \text{Ker } \mathbf{D}^{(k)}$ and therefore $\mathbf{v} \in \bigcap_{k=0}^{\infty} \text{Ker } \mathbf{D}^{(k)}$. So \mathbf{v} can only be the null vector, which is absurd since $\|\mathbf{v}\| = 1$. We deduce from this contradiction that the limit of the sequence of integers $(d(n))_{n \in \mathbb{N}}$ is 0. Therefore there exists a nonnegative integer N such that $d(N) = 0$. □

Appendix 3.B Maclaurin series

The lemmas in this subsection deal with the following setting.

Let m be a positive integer and let \mathbf{M} be a continuous mapping from \mathbb{R} to \mathcal{M}_m , the set of $m \times m$ matrices. Assume \mathbf{M} admits the following Maclaurin series:

$$\mathbf{M}(t) = \sum_{k=0}^N a_k(t) \mathbf{A}_k + \mathbf{B}(t). \quad (3.34)$$

In the above expression, N is a nonnegative integer and for every $k \in \llbracket 0, N \rrbracket$:

- (a) a_k is a continuous mapping $\mathbb{R} \rightarrow \mathbb{R}$ such that for all $t \neq 0$, $a_k(t) \neq 0$;
- (b) for every nonnegative integer $l < k$, $a_k(t) = o(|a_l(t)|)$ when $t \rightarrow 0$;
- (c) \mathbf{A}_k is a non-null symmetric $m \times m$ matrix.

\mathbf{B} is a continuous mapping $\mathbb{R} \rightarrow \mathcal{M}_m$ such that for every $t \in \mathbb{R}$, $\mathbf{B}(t)$ is a symmetric matrix and when $t \rightarrow 0$, $\|\mathbf{B}(t)\| = o(|a_N(t)|)$.

Lemma 3.14. *Consider (3.34). If $\cap_{k=0}^N \text{Ker } \mathbf{A}_k$ is the trivial vector space and if there exists $T > 0$ such that for all $t \in (-T, T)$ $\mathbf{M}(t)$ is nonsingular, then when $t \rightarrow 0$, $\|\mathbf{M}(t)^{-1}\| = O(|a_N(t)|^{-1})$.*

Proof. Assume that $\cap_{k=0}^N \text{Ker } \mathbf{A}_k$ is the trivial vector space and that there exists $T > 0$ such that for all $t \in (-T, T)$, $\mathbf{M}(t)$ is a nonsingular matrix.

If $N = 0$, then \mathbf{A}_0 is nonsingular and the conclusion is trivial.

If $N \geq 1$, we may assume without loss of generality that $\cap_{k=0}^{N-1} \text{Ker } \mathbf{A}_k$ is a nontrivial vector space, otherwise we could replace N by $N - 1$ and $\mathbf{B}(t)$ by $\{a_N(t)\mathbf{A}_N + \mathbf{B}(t)\}$ for all $t \in \mathbb{R}$.

Let d_N be the codimension of $\cap_{k=0}^{N-1} \text{Ker } \mathbf{A}_k$. Let \mathbf{W}_N be an $m \times (m - d_N)$ matrix whose columns form an orthonormal basis of $\cap_{k=0}^{N-1} \text{Ker } \mathbf{A}_k$, and let \mathbf{P}_N be an $m \times d_N$ matrix whose columns form an orthonormal basis of its orthogonal complement. Then $(\mathbf{P}_N \mathbf{W}_N)$ is an orthogonal matrix. For all $t \in \mathbb{R}$, let us replace $\mathbf{M}(t)$ by $(\mathbf{P}_N \mathbf{W}_N)^\top \mathbf{M}(t) (\mathbf{P}_N \mathbf{W}_N)$. Because $(\mathbf{P}_N \mathbf{W}_N)$ is an orthogonal matrix, the Frobenius norm of $\mathbf{M}(t)^{-1}$ is unchanged. Naturally, for all $k \in \llbracket 0, N \rrbracket$, \mathbf{A}_k is replaced by $(\mathbf{P}_N \mathbf{W}_N)^\top \mathbf{A}_k (\mathbf{P}_N \mathbf{W}_N)$ and for every $t \in \mathbb{R}$, $\mathbf{B}(t)$ is replaced by $(\mathbf{P}_N \mathbf{W}_N)^\top \mathbf{B}(t) (\mathbf{P}_N \mathbf{W}_N)$.

Now, for every $k \in \llbracket 1, N \rrbracket$, \mathbf{A}_k can be decomposed into blocks – a $d_N \times d_N$ block \mathbf{A}'_k , an $(m - d_N) \times (m - d_N)$ block \mathbf{A}''_k and a $d_N \times (m - d_N)$ block \mathbf{A}'''_k :

$$\mathbf{A}_k = \begin{pmatrix} \mathbf{A}'_k & \mathbf{A}'''_k \\ (\mathbf{A}'''_k)^\top & \mathbf{A}''_k \end{pmatrix} \quad (3.35)$$

For all $t \in \mathbb{R}$, $\mathbf{B}(t)$ can be decomposed in a similar manner (here the $'$ notation is used to distinguish the blocks, not to express some derivative with respect to t):

$$\mathbf{B}(t) = \begin{pmatrix} \mathbf{B}(t)' & \mathbf{B}(t)''' \\ (\mathbf{B}(t)''')^\top & \mathbf{B}(t)'' \end{pmatrix} \quad (3.36)$$

Now, for any symmetric nonsingular matrix

$$\mathbf{C} = \begin{pmatrix} \mathbf{C}' & \mathbf{C}''' \\ (\mathbf{C}''')^\top & \mathbf{C}'' \end{pmatrix}, \quad (3.37)$$

denoting by $\mathbf{S} := \left\{ \mathbf{C}' - \mathbf{C}''' (\mathbf{C}'')^{-1} (\mathbf{C}''')^\top \right\}$ the Schur complement of \mathbf{C}'' , the inverse of \mathbf{C} is

$$\mathbf{C}^{-1} = \begin{pmatrix} \mathbf{I} & \mathbf{0} \\ -(\mathbf{C}'')^{-1} (\mathbf{C}''')^\top & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{S}^{-1} & \mathbf{0} \\ \mathbf{0} & (\mathbf{C}'')^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{I} & -\mathbf{C}''' (\mathbf{C}'')^{-1} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}. \quad (3.38)$$

For every $k \in \llbracket 0, N-1 \rrbracket$, \mathbf{A}_k'' and \mathbf{A}_k''' are null. For all $t \in (-T, T)$, $\mathbf{M}(t)$ is nonsingular. Its lower $(m - d_N) \times (m - d_N)$ block is $\{a_N(t)\mathbf{A}_N'' + \mathbf{B}(t)''\}$ and its Schur complement $\mathbf{S}_N(t)$ is

$$\begin{aligned} \mathbf{S}_N(t) &:= -\{a_N(t)\mathbf{A}_N''' + \mathbf{B}(t)'''\} \{a_N(t)\mathbf{A}_N'' + \mathbf{B}(t)''\}^{-1} \{a_N(t)\mathbf{A}_N''' + \mathbf{B}(t)'''\}^\top \\ &\quad + \left\{ \sum_{k=0}^N a_k(t)\mathbf{A}_k' + \mathbf{B}(t)' \right\}. \end{aligned} \quad (3.39)$$

Because we are dealing with the finite dimensional vector space of matrices of size $m \times m$, all norms are equivalent. In particular, the Frobenius norm is equivalent to the algebra norm

$$\mathbf{A} \mapsto \sup \left\{ \sqrt{\boldsymbol{\xi}^\top \mathbf{A} \boldsymbol{\xi} / \boldsymbol{\xi}^\top \boldsymbol{\xi}} : \boldsymbol{\xi} \in \mathbb{R}^m \setminus \{\mathbf{0}_m\} \right\}.$$

So there exists a constant $C_m \in (0, +\infty)$ such that for every $t \in (-T, T)$,

$$\begin{aligned} \|\mathbf{M}(t)^{-1}\| &\leq C_m \left(\|\mathbf{I}_m\| + \left\| \{a_N(t)\mathbf{A}_N''' + \mathbf{B}(t)'''\} \{a_N(t)\mathbf{A}_N'' + \mathbf{B}(t)''\}^{-1} \right\| \right)^2 \\ &\quad \left(\|\mathbf{S}_N(t)^{-1}\| + \left\| \{a_N(t)\mathbf{A}_N'' + \mathbf{B}(t)''\}^{-1} \right\| \right). \end{aligned} \quad (3.40)$$

\mathbf{A}_N'' is nonsingular, otherwise $\cap_{k=0}^N \text{Ker } \mathbf{A}_k$ would be nontrivial. This means that the norm of the matrix $\{a_N(t)\mathbf{A}_N''' + \mathbf{B}(t)'''\} \{a_N(t)\mathbf{A}_N'' + \mathbf{B}(t)''\}^{-1}$ is bounded when $t \rightarrow 0$. Because of Equation (3.40), this implies that there exists $T_N > 0$ and $\lambda_N > 0$ such that for all $t \in (-T_N, T_N)$,

$$\lambda_N \|\mathbf{M}(t)^{-1}\| \leq |a_N(t)|^{-1} + \|\mathbf{S}_N(t)^{-1}\|. \quad (3.41)$$

Our goal is to use Equation (3.41) recursively, by having $\mathbf{S}_N(t)$ take the place of $\mathbf{M}(t)$. To achieve this, a new expression of $\mathbf{S}_N(t)$ is required.

$$\mathbf{S}_N(t) = \sum_{k=0}^{N-1} a_k(t)\mathbf{A}_k' + \mathbf{B}_N(t), \quad (3.42)$$

where

$$\begin{aligned} \mathbf{B}_N(t) &:= -\{a_N(t)\mathbf{A}_N''' + \mathbf{B}(t)'''\} \{a_N(t)\mathbf{A}_N'' + \mathbf{B}(t)''\}^{-1} \{a_N(t)\mathbf{A}_N''' + \mathbf{B}(t)'''\}^\top \\ &\quad + a_N(t)\mathbf{A}_N' + \mathbf{B}(t)'. \end{aligned} \quad (3.43)$$

It turns out that when $t \rightarrow 0$, the norm of $\mathbf{B}_N(t)$ is $O(|a_N(t)|)$. This is due to the fact mentioned above that $\{a_N(t)\mathbf{A}_N''' + \mathbf{B}(t)'''\} \{a_N(t)\mathbf{A}_N'' + \mathbf{B}(t)''\}^{-1}$ is bounded when $t \rightarrow 0$.

Furthermore, $\cap_{k=0}^{N-1} \text{Ker } \mathbf{A}_k'$ is the trivial vector space. Indeed, let $\mathbf{v}_1 \in \cap_{k=0}^{N-1} \text{Ker } \mathbf{A}_k'$. Then for any vector $\mathbf{v}_2 \in \mathbb{R}^{m-d_N}$, $(\mathbf{v}_1, \mathbf{v}_2)^\top \cap_{k=0}^{N-1} \text{Ker } \mathbf{A}_k$. Independently from this, for any vector $\mathbf{v}_3 \in \mathbb{R}^{d_N}$, $(\mathbf{v}_3, \mathbf{0}_{m-d_N})^\top$ belongs to the orthogonal complement of $\cap_{k=0}^{N-1} \text{Ker } \mathbf{A}_k$. So

$(\mathbf{v}_1, \mathbf{0}_{m-d_N})^\top$ belongs both to $\cap_{k=0}^{N-1} \text{Ker } \mathbf{A}_k$ and its orthogonal complement: it is the null vector. Therefore $\mathbf{v}_1 = \mathbf{0}_{d_N}$.

The two paragraphs above show that Equation (3.42) is formally similar to Equation (3.34): the role of $\mathbf{M}(t)$ is held by $\mathbf{S}_N(t)$, the role of N by $N - 1$, the role of the \mathbf{A}_k s by the \mathbf{A}'_k s and the role of $\mathbf{B}(t)$ by $\mathbf{B}_N(t)$.

Therefore an equation similar to (3.41) can be derived: there exist $T_{N-1} > 0$ and $\lambda_{N-1} > 0$ such that for all $t \in (-T_{N-1}, T_{N-1})$,

$$\lambda_{N-1} \|\mathbf{S}_N(t)^{-1}\| \leq |a_{N-1}(t)|^{-1} + \|\mathbf{S}_{N-1}(t)^{-1}\|. \quad (3.44)$$

Here, $\mathbf{S}_{N-1}(t)$ is defined with respect to $\mathbf{S}_N(t)$ the same way $\mathbf{S}_N(t)$ was defined with respect to $\mathbf{M}(t)$.

Recursive application of this reasoning until 0 is reached yields the result. \square

Lemma 3.15. *Consider (3.34) with $N = 1$. If $\text{Ker } \mathbf{A}_0 \cap \text{Ker } \mathbf{A}_1$ is the trivial vector space, then there exists $T > 0$ such that for all $t \in (-T, T)$, $\mathbf{M}(t)$ is nonsingular.*

Proof. We use the same notations as in the proof of Lemma 3.14 and redefine matrices the same way: $\mathbf{M}(t) := (\mathbf{P}_1 \mathbf{W}_1)^\top \mathbf{M}(t) (\mathbf{P}_1 \mathbf{W}_1)$, $\mathbf{B}(t) := (\mathbf{P}_1 \mathbf{W}_1)^\top \mathbf{B}(t) (\mathbf{P}_1 \mathbf{W}_1)$, $\mathbf{A}_0 := (\mathbf{P}_1 \mathbf{W}_1)^\top \mathbf{A}_0 (\mathbf{P}_1 \mathbf{W}_1)$, $\mathbf{A}_1 := (\mathbf{P}_1 \mathbf{W}_1)^\top \mathbf{A}_1 (\mathbf{P}_1 \mathbf{W}_1)$. For all $t \in \mathbb{R}$, the determinant of $\mathbf{M}(t)$ is the product of the determinants of $\mathbf{S}_1(t)$ and of $\{a_1(t)\mathbf{A}_1'' + \mathbf{B}(t)''\}$. Because $\text{Ker } \mathbf{A}_0 \cap \text{Ker } \mathbf{A}_1$ is the trivial vector space, \mathbf{A}_1'' is nonsingular. When $|t|$ is small enough therefore, $\{a_1(t)\mathbf{A}_1'' + \mathbf{B}(t)''\}$ is nonsingular too. Moreover, $\lim_{t \rightarrow 0} a_0^{-1}(t)\mathbf{S}_1(t) = \mathbf{A}'_0$ is also nonsingular. These two facts imply that when $|t|$ is small enough, the determinant of $\mathbf{M}(t)$ is non-null and $\mathbf{M}(t)$ is nonsingular. \square

Lemma 3.16. *Consider (3.34). If $\cap_{k=0}^N \text{Ker } \mathbf{A}_k$ is the trivial vector space, if the vector space $\cap_{k=0}^{N-1} \text{Ker } \mathbf{A}_k$ is non-trivial, and if there exists $T > 0$ such that for all $t \in (-T, T)$, $\mathbf{M}(t)$ is positive definite, then there exists a hyperplane \mathcal{H} of \mathbb{R}^m such that for all $\mathbf{v} \in \mathbb{R}^m \setminus \mathcal{H}$,*

$$\liminf_{t \rightarrow 0} \mathbf{v} \mathbf{M}(t)^{-1} \mathbf{v} / \|\mathbf{M}(t)^{-1}\| > 0. \quad (3.45)$$

Proof. This result is trivial if $N = 0$ ($\cap_{k=0}^{-1} \text{Ker } \mathbf{A}_k$ is an intersection over an empty set, so we take it by convention to be \mathbb{R}^m). If $N \geq 1$, it follows from the proof of Lemma 3.14. Indeed, the requirements of this lemma are stronger than those of Lemma 3.14, so all intermediate results of its proof are valid. Consider Equation (3.38) while assuming \mathbf{C} is positive definite. In the right member, the matrices on the left and on the right are the transpose of one another, so the middle matrix is necessarily positive definite. In particular, both \mathbf{S}^{-1} and $(\mathbf{C}'')^{-1}$ are positive definite. Any vector $\mathbf{v} \in \mathbb{R}^m$ can be decomposed as $\mathbf{v} = (\mathbf{v}', \mathbf{v}'')^\top$ with $\mathbf{v}' \in \mathbb{R}^{d_N}$ and $\mathbf{v}'' \in \mathbb{R}^{m-d_N}$. This decomposition yields a lower bound: $\mathbf{v}^\top \mathbf{C}^{-1} \mathbf{v} \geq (\mathbf{v}'')^\top (\mathbf{C}'')^{-1} \mathbf{v}''$. Here, \mathbf{C} is $\mathbf{M}(t)$, \mathbf{S} is $\mathbf{S}_N(t)$ and \mathbf{C}'' is $a_N(t)\mathbf{A}_N'' + \mathbf{B}''(t)$. Let us recall that \mathbf{A}_N'' is nonsingular and $\|\mathbf{B}''(t)\| = o(|a_N(t)|)$ when $t \rightarrow 0$. So as long as \mathbf{v} is not orthogonal to $\cap_{k=0}^{N-1} \text{Ker } \mathbf{A}_k$, \mathbf{v}'' is non-zero and there exists $\tilde{\lambda}_N(\mathbf{v}) > 0$ such that when $|t|$ is small enough, $\mathbf{v}^\top \mathbf{M}(t)^{-1} \mathbf{v} \geq \tilde{\lambda}_N(\mathbf{v}) |a_N(t)|^{-1}$. Then Lemma 3.14 yields the result for any hyperplane \mathcal{H} of \mathbb{R}^m that contains the orthogonal complement of $\cap_{k=0}^{N-1} \text{Ker } \mathbf{A}_k$. \square

Lemma 3.17. *If $\bigcap_{k=0}^N \text{Ker } \mathbf{A}_k \neq \text{Ker } \mathbf{A}_0$, if $\bigcap_{k=0}^{N-1} \text{Ker } \mathbf{A}_k = \text{Ker } \mathbf{A}_0$, and if there exists $T > 0$ such that for all $t \in (-T, T)$, $\mathbf{M}(t)$ is positive definite, then the largest eigenvalue $v_1(t)$ and the second largest eigenvalue $v_2(t)$ of $\mathbf{M}(t)$ have the following behavior when $t \rightarrow 0$:*

$$(a) \ v_1(t)^{-1} = O(|a_0(t)|^{-1});$$

$$(b) \ v_2(t)^{-1} = O(|a_N(t)|^{-1}).$$

Proof. It is equivalent to prove that in this situation, there exists $\lambda > 0$ such that when $|t|$ is sufficiently small $v_1(t) \geq \lambda|a_0(t)|$ and $v_2(t) \geq \lambda|a_N(t)|$.

When $t \rightarrow 0$, we have $a_0(t)^{-1}\mathbf{M}(t) \rightarrow \mathbf{A}_0$, so \mathbf{A}_0 is either positive or negative semi-definite. Since a_0 is continuous and non-null everywhere except possibly at 0, its sign is therefore constant: nonnegative if \mathbf{A}_0 is positive semi-definite and nonpositive if \mathbf{A}_0 is negative semi-definite. Without loss of generality, let us assume that \mathbf{A}_0 is positive semi-definite and that a_0 is nonnegative. $a_0(t)^{-1}\mathbf{M}(t) \rightarrow \mathbf{A}_0$ implies that $a_0(t)^{-1}v_1(t)$ converges to \mathbf{A}_0 's greatest eigenvalue, which is strictly greater than 0 because \mathbf{A}_0 is non-null. This implies the first result.

Now, since \mathbf{A}_0 is non-null, its rank is greater or equal to 1. If it is greater or equal to 2, then $a_0(t)^{-1}v_2(t)$ converges to the second greatest eigenvalue of \mathbf{A}_0 , so $v_2^{-1}(t) = O(a_0(t)^{-1})$ and the second result holds a fortiori.

Assume from now on that \mathbf{A}_0 has rank 1. For every nonnegative integer $k < N$, \mathbf{A}_k shares \mathbf{A}_0 's kernel, so \mathbf{A}_k is proportional to \mathbf{A}_0 . We may therefore assume without loss of generality that $N = 1$. Since \mathbf{A}_0 is a symmetric positive semi-definite matrix with rank 1, there exists a vector \mathbf{a}_0 such that $\mathbf{A}_0 = \mathbf{a}_0\mathbf{a}_0^\top$.

Choose for all $t \in (-T, T)$ a unit eigenvector $\mathbf{V}_1(t)$ corresponding to the eigenvalue $v_1(t)$ of $\mathbf{M}(t)$. Then choose a unit eigenvector $\mathbf{V}_2(t)$ corresponding to the eigenvalue $v_2(t)$ such that $\mathbf{V}_1(t)^\top \mathbf{V}_2(t) = 0$ (it is always possible to choose $\mathbf{V}_2(t)$ that way because $\mathbf{M}(t)$ is symmetric). When $t \rightarrow 0$, $\mathbf{V}_1(t) \rightarrow \mathbf{a}_0/\|\mathbf{a}_0\|$. Since $\mathbf{V}_1(t)^\top \mathbf{V}_2(t) = 0$ for all $t \in (-T, T)$, we have $\lim_{t \rightarrow 0} \mathbf{a}_0^\top \mathbf{V}_2(t) = 0$.

$$\begin{aligned} v_2(t) &= a_0(t) (\mathbf{a}_0^\top \mathbf{V}_2(t))^2 + a_1(t) \mathbf{V}_2(t)^\top \mathbf{A}_1 \mathbf{V}_2(t) + \mathbf{V}_2(t)^\top \mathbf{B}(t) \mathbf{V}_2(t) \\ &\geq a_1(t) \{ \mathbf{V}_2(t)^\top \mathbf{A}_1 \mathbf{V}_2(t) + o(1) \}. \end{aligned} \quad (3.46)$$

Because $\mathbf{M}(t)$ is positive definite for all $t \in (-T, T)$, the restriction of \mathbf{A}_1 to $\text{Ker } \mathbf{A}_0$ is either positive semi-definite (making a_1 nonnegative) or negative semi-definite (making a_1 nonpositive). Moreover, it is non-null.

Since $v_2(t) = \max\{\boldsymbol{\xi}^\top \mathbf{M}(t) \boldsymbol{\xi} \mid \boldsymbol{\xi} \in \mathbb{R}^m \text{ and } \|\boldsymbol{\xi}\| = 1 \text{ and } \boldsymbol{\xi}^\top \mathbf{v}_1(t) = 0\}$, the above implies the following: $\liminf_{t \rightarrow 0} |a_1(t)|^{-1} v_2(t) > 0$. So the second result also holds when the rank of \mathbf{A}_0 is 1.

□

Appendix 3.C Spectral decomposition

For the following lemmas, we need to set up a few notations. First, denote by \widehat{K}_r the r -dimensional Fourier transform of the isotropic correlation kernel K :

$$\widehat{K}_r(\boldsymbol{\omega}) = (2\pi)^{-r} \int_{\mathbb{R}^r} K(\|\boldsymbol{x}\|) e^{-i\langle \boldsymbol{\omega}, \boldsymbol{x} \rangle} d\boldsymbol{x} \quad \text{and} \quad K(\|\boldsymbol{x}\|) = \int_{\mathbb{R}^r} \widehat{K}_r(\boldsymbol{\omega}) e^{i\langle \boldsymbol{\omega}, \boldsymbol{x} \rangle} d\boldsymbol{\omega}. \quad (3.47)$$

For all $\theta \in (0, +\infty)$, using the correlation kernel $K_\theta(\cdot) = K(\cdot/\theta)$, the correlation matrix $\boldsymbol{\Sigma}_\theta$ is such that $\forall \boldsymbol{\xi} \in \mathbb{R}^n$:

$$\boldsymbol{\xi}^\top \boldsymbol{\Sigma}_\theta \boldsymbol{\xi} = \sum_{j,k=1}^n \xi_j \xi_k K\left(\frac{\|\boldsymbol{x}^{(j)} - \boldsymbol{x}^{(k)}\|}{\theta}\right) = \int_{\mathbb{R}^r} \widehat{K}_r(\boldsymbol{\omega}) \left| \sum_{j=1}^n \xi_j e^{i\langle \boldsymbol{\omega}, \frac{\boldsymbol{x}^{(j)}}{\theta} \rangle} \right|^2 d\boldsymbol{\omega} = M_r \theta^r I_\theta(\boldsymbol{\xi}). \quad (3.48)$$

The factors in the last equality depend on the kernel and are given in Table 3.4.

Kernel	M_r	$I_\theta(\boldsymbol{\xi})$
Matérn	$\frac{\Gamma(\nu + \frac{r}{2})(2\sqrt{\nu})^{2\nu}}{\pi^{\frac{r}{2}} \Gamma(\nu)}$	$\int_{\mathbb{R}^r} (4\nu + \theta^2 \ \boldsymbol{s}\ ^2)^{-\frac{r}{2} - \nu} \left \sum_{j=1}^n \xi_j e^{i\langle \boldsymbol{s}, \boldsymbol{x}^{(j)} \rangle} \right ^2 d\boldsymbol{s}$
Rational Quadratic	$\frac{2^{1-\nu}}{(2\pi)^{\frac{r}{2}} \Gamma(\nu)}$	$\int_{\mathbb{R}^r} (\theta \ \boldsymbol{s}\)^{\nu - \frac{r}{2}} \mathcal{K}_{\nu - \frac{r}{2}}(\theta \ \boldsymbol{s}\) \left \sum_{j=1}^n \xi_j e^{i\langle \boldsymbol{s}, \boldsymbol{x}^{(j)} \rangle} \right ^2 d\boldsymbol{s}$
Squared Exponential	$(2\sqrt{\pi})^{-r}$	$\int_{\mathbb{R}^r} \exp\left(-\frac{\theta^2 \ \boldsymbol{s}\ ^2}{4}\right) \left \sum_{j=1}^n \xi_j e^{i\langle \boldsymbol{s}, \boldsymbol{x}^{(j)} \rangle} \right ^2 d\boldsymbol{s}$

Table 3.4 – M_r and $I_\theta(\boldsymbol{\xi})$ for the three considered correlation kernel families. \mathcal{K}_ν is the modified Bessel function of second kind with parameter ν [Abramowitz and Stegun, 1964] (9.6.)

The spectral decomposition of correlation kernels is a powerful tool.

To use it, recall Proposition 1.31:

Lemma 3.18. *Let μ be a positive measure on \mathbb{R}^r with finite non-null total mass that is absolutely continuous with respect to the Lebesgue measure. Then the mapping $K : \mathbb{R}^r \rightarrow \mathbb{R}$ defined by*

$$K(\boldsymbol{x}) = \int_{\mathbb{R}^r} e^{i\langle \boldsymbol{\omega}, \boldsymbol{x} \rangle} d\mu(\boldsymbol{\omega}) \quad (3.49)$$

is positive definite. Moreover, for any $\boldsymbol{\xi} \in \mathbb{R}^n \setminus \{\mathbf{0}_n\}$,

$$\sum_{k,l \in [1,n]} \xi_k \xi_l K(\boldsymbol{x}^{(k)} - \boldsymbol{x}^{(l)}) > 0. \quad (3.50)$$

Let us use spectral decomposition to show this useful fact about Matérn kernels:

Lemma 3.19. *For Matérn kernels, when $\theta \rightarrow +\infty$ $\|\boldsymbol{\Sigma}_\theta^{-1}\| = O(\theta^{2\nu})$.*

Proof. When $\theta \geq 2\sqrt{\nu}$, for any $\boldsymbol{\xi} \in \mathbb{R}^n$,

$$I_\theta(\boldsymbol{\xi}) \geq 2^{-\frac{r}{2} - \nu} \theta^{-r-2\nu} \int_{\|\boldsymbol{s}\| \geq 1} \|\boldsymbol{s}\|^{-r-2\nu} \left| \sum_{j=1}^n \xi_j e^{i\langle \boldsymbol{s}, \boldsymbol{x}^{(j)} \rangle} \right|^2 d\boldsymbol{s}. \quad (3.51)$$

Define the mapping $K^{aux} : \mathbb{R}^r \rightarrow \mathbb{R}$ by

$$K^{aux}(\mathbf{x}) = \int_{\mathbb{R}^r} e^{i\langle \boldsymbol{\omega} | \mathbf{x} \rangle} \|\mathbf{s}\|^{-r-2\nu} \mathbf{1}_{\|\mathbf{s}\| \geq 1} d\mathbf{s}. \quad (3.52)$$

By Lemma 3.18, the $n \times n$ matrix \mathbf{M} with (i, i') -th element $K^{aux}(\mathbf{x}^{(i)} - \mathbf{x}^{(i')})$ is positive definite. For any $\boldsymbol{\xi} \in \mathbb{R}^n$, $\int_{\|\mathbf{s}\| \geq 1} \|\mathbf{s}\|^{-r-2\nu} \left| \sum_{j=1}^n \xi_j e^{i\langle \mathbf{s} | \mathbf{x}^{(j)} \rangle} \right|^2 d\mathbf{s} = \boldsymbol{\xi}^\top \mathbf{M} \boldsymbol{\xi}$. Denote by M the smallest eigenvalue of \mathbf{M} . For any $\boldsymbol{\xi} \in \mathbb{R}^n$, when $\theta \geq 2\sqrt{\nu}$, $I_\theta(\boldsymbol{\xi}) \geq 2^{-\frac{r}{2}-\nu} M \|\boldsymbol{\xi}\|^2 \theta^{-r-2\nu}$. Equation (3.48) implies the result. \square

More generally, it can be used to study the behavior of the reference prior. From Equation (3.48), we obtain that $\forall \theta \in (0, +\infty)$, $\forall \boldsymbol{\xi} \in \mathbb{R}^n$:

$$\boldsymbol{\xi}^\top \left(\frac{d}{d\theta} \boldsymbol{\Sigma}_\theta \right) \boldsymbol{\xi} = M_r r \theta^{r-1} I_\theta(\boldsymbol{\xi}) + M_r \theta^r \frac{d}{d\theta} I_\theta(\boldsymbol{\xi}). \quad (3.53)$$

The next three lemmas are used to prove the second assertion of Proposition 3.7. Since the proof varies for each of the three different kernel families considered, I_θ is written I_θ^M for Matérn kernels, I_θ^{RQ} for Rational Quadratic kernels and I_θ^{SE} for the Squared Exponential kernel.

Lemma 3.20. *For Matérn kernels, $\mathbf{F}_\theta := r\theta^{-1}\boldsymbol{\Sigma}_\theta - \frac{d}{d\theta}\boldsymbol{\Sigma}_\theta$ is a symmetric positive definite matrix. Furthermore, for any $\boldsymbol{\xi} \in \mathbb{R}^n$, $\boldsymbol{\xi}^\top \mathbf{F}_\theta \boldsymbol{\xi} \leq (2\nu + r)\theta^{-1}\boldsymbol{\xi}^\top \boldsymbol{\Sigma}_\theta \boldsymbol{\xi}$.*

Proof. For any $\theta \in (0, +\infty)$ and any $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n)^\top \in \mathbb{R}^n$,

$$\begin{aligned} \frac{d}{d\theta} I_\theta^M(\boldsymbol{\xi}) &= (-2) \left(\frac{r}{2} + \nu \right) \theta \int_{\mathbb{R}^r} \|\mathbf{s}\|^2 (4\nu + \theta^2 \|\mathbf{s}\|^2)^{-\frac{r}{2}-\nu-1} \left| \sum_{j=1}^n \xi_j e^{i\langle \mathbf{s} | \mathbf{x}^{(j)} \rangle} \right|^2 d\mathbf{s} \\ &= -(2\nu + r)\theta^{-1} \int_{\mathbb{R}^r} \frac{\theta^2 \|\mathbf{s}\|^2}{4\nu + \theta^2 \|\mathbf{s}\|^2} (4\nu + \theta^2 \|\mathbf{s}\|^2)^{-\frac{r}{2}-\nu} \left| \sum_{j=1}^n \xi_j e^{i\langle \mathbf{s} | \mathbf{x}^{(j)} \rangle} \right|^2 d\mathbf{s} \end{aligned} \quad (3.54)$$

Since the ratio in the integrand is smaller than 1, for any $\theta \in (0, +\infty)$ and any non-null vector $\boldsymbol{\xi} \in \mathbb{R}^n$,

$$0 < -\frac{d}{d\theta} I_\theta^M(\boldsymbol{\xi}) \leq (2\nu + r)\theta^{-1} I_\theta^M(\boldsymbol{\xi}) \quad (3.55)$$

Combining Equations (3.48), (3.53) and (3.55) yields the result. \square

Lemma 3.21. *For Rational Quadratic isotropic correlation kernels, $\mathbf{F}_\theta := r\theta^{-1}\boldsymbol{\Sigma}_\theta - \frac{d}{d\theta}\boldsymbol{\Sigma}_\theta$ is a symmetric positive definite matrix. If θ is large enough, it verifies $\forall \boldsymbol{\xi} \in \mathbb{R}^n$, $\boldsymbol{\xi}^\top \mathbf{F}_\theta \boldsymbol{\xi} \leq (r+2)\boldsymbol{\xi}^\top \boldsymbol{\Sigma}_\theta \boldsymbol{\xi}$.*

Proof. In the following, denote by \mathcal{K}_ν the modified Bessel function of second kind with parameter ν . [Abramowitz and Stegun, 1964](9.6.)

For any $\theta \in (0, +\infty)$ and any $\boldsymbol{\xi} \in \mathbb{R}^n$,

$$\frac{d}{d\theta} I_\theta^{RQ}(\boldsymbol{\xi}) = \int_{\mathbb{R}^r} \|\mathbf{s}\| \frac{d}{dz} \left\{ z^{\nu-\frac{r}{2}} \mathcal{K}_{\nu-\frac{r}{2}}(z) \right\}_{z=\theta\|\mathbf{s}\|} \left| \sum_{j=1}^n \xi_j e^{i\langle \mathbf{s} | \mathbf{x}^{(j)} \rangle} \right|^2 d\mathbf{s}. \quad (3.56)$$

We now compute $\frac{d}{dz} \left\{ z^{\nu-\frac{r}{2}} \mathcal{K}_{\nu-\frac{r}{2}}(z) \right\}$. Following Abramowitz and Stegun [1964] (9.6.28.),

$$\frac{d}{dz} \left\{ z^{\nu-\frac{r}{2}} \mathcal{K}_{\nu-\frac{r}{2}}(z) \right\} = \begin{cases} -z^{\nu-\frac{r}{2}} \mathcal{K}_{\nu-\frac{r}{2}-1}(z) & \text{if } \nu - \frac{r}{2} \geq 0. \\ -z^{\nu-\frac{r}{2}} \mathcal{K}_{\frac{r}{2}-\nu-1}(z) + (2\nu-r)z^{\nu-\frac{r}{2}-1} \mathcal{K}_{\frac{r}{2}-\nu}(z) & \text{if } \nu - \frac{r}{2} < 0. \end{cases} \quad (3.57)$$

Combining this with Equations (3.48) and (3.53) proves that \mathbf{F}_θ is positive definite.

We now have to deal with the behavior of $\mathcal{K}_{|\nu-\frac{r}{2}|-1}(z)$ when $z \rightarrow 0$ and when $z \rightarrow +\infty$.

Let us start with $z \rightarrow 0$. Using Abramowitz and Stegun [1964] (9.6.9.), we obtain:

$$\mathcal{K}_{|\nu-\frac{r}{2}|-1}(z) \sim \begin{cases} \frac{z}{2(|\nu-\frac{r}{2}|-1)} \mathcal{K}_{\nu-\frac{r}{2}}(z) & \text{if } |\nu - \frac{r}{2}| > 1. \\ -z \log(z) \mathcal{K}_{\nu-\frac{r}{2}}(z) & \text{if } |\nu - \frac{r}{2}| = 1. \\ \frac{\Gamma(1-|\nu-\frac{r}{2}|)}{2^{|\nu-\frac{r}{2}|-1} \Gamma(|\nu-\frac{r}{2}|)} z^{2\nu-r-1} \mathcal{K}_{\nu-\frac{r}{2}}(z) & \text{if } 0 < |\nu - \frac{r}{2}| < 1. \\ -\frac{1}{z \log(z)} \mathcal{K}_{\nu-\frac{r}{2}}(z) & \text{if } |\nu - \frac{r}{2}| = 0. \end{cases} \quad (3.58)$$

So, for any $\nu > 0$, there exists $a_{r,\nu} > 0$ such that, as long as z is small enough, $\mathcal{K}_{|\nu-\frac{r}{2}|-1}(z) \leq a_{r,\nu} z^{-1} \mathcal{K}_{\nu-\frac{r}{2}}(z)$.

Moreover, Abramowitz and Stegun [1964] (9.7.2.) states that when $z \rightarrow +\infty$, $\mathcal{K}_{|\nu-\frac{r}{2}|-1}(z) \sim \mathcal{K}_{\nu-\frac{r}{2}}(z) \sim \exp(-z) \sqrt{\pi} / \sqrt{2z}$.

Because $\mathcal{K}_{|\nu-\frac{r}{2}|-1}$ is a continuous function on $(0, +\infty)$, the two results above imply that

$$\forall \lambda > 1 \quad \exists a'_{r,\nu} > 0 \quad \forall z > 0 \quad z \mathcal{K}_{|\nu-\frac{r}{2}|-1} \leq \max(a'_{r,\nu}, \lambda z) \mathcal{K}_{\nu-\frac{r}{2}}(z). \quad (3.59)$$

Now, $\frac{d}{d\theta} I_\theta^{RQ}(\boldsymbol{\xi}) = J_\theta^1(\boldsymbol{\xi}) + J_\theta^2(\boldsymbol{\xi})$, with

$$\begin{aligned} J_\theta^1(\boldsymbol{\xi}) &:= \int_{\|\mathbf{s}\| \leq 1} \|\mathbf{s}\| \frac{d}{dz} \left\{ z^{\nu-\frac{r}{2}} \mathcal{K}_{\nu-\frac{r}{2}}(z) \right\}_{z=\theta\|\mathbf{s}\|} \left| \sum_{j=1}^n \xi_j e^{i\langle \mathbf{s} | \mathbf{x}^{(j)} \rangle} \right|^2 d\mathbf{s}; \\ J_\theta^2(\boldsymbol{\xi}) &:= \int_{\|\mathbf{s}\| > 1} \|\mathbf{s}\| \frac{d}{dz} \left\{ z^{\nu-\frac{r}{2}} \mathcal{K}_{\nu-\frac{r}{2}}(z) \right\}_{z=\theta\|\mathbf{s}\|} \left| \sum_{j=1}^n \xi_j e^{i\langle \mathbf{s} | \mathbf{x}^{(j)} \rangle} \right|^2 d\mathbf{s}. \end{aligned} \quad (3.60)$$

Set $0 < \epsilon < 1$. When $z \in [\epsilon, 1]$, $\frac{d}{dz} \left\{ z^{\nu-\frac{r}{2}} \mathcal{K}_{\nu-\frac{r}{2}}(z) \right\}$ is bounded away from 0, so there exists m_ϵ such that

$$\theta |J_\theta^1(\boldsymbol{\xi})| \geq m_\epsilon \int_{\frac{\epsilon}{\theta} < \|\mathbf{s}\| < \frac{1}{\theta}} \left| \sum_{j=1}^n \xi_j e^{i\langle \mathbf{s} | \mathbf{x}^{(j)} \rangle} \right|^2 d\mathbf{s} = m_\epsilon \theta^{-r} \int_{\mathbf{1}_{\epsilon < \|\boldsymbol{\omega}\| < 1}} \left| \sum_{j=1}^n \xi_j e^{i\langle \boldsymbol{\omega} | \frac{\mathbf{x}^{(j)}}{\theta} \rangle} \right|^2 d\boldsymbol{\omega}. \quad (3.61)$$

The Lebesgue measure on $\{\boldsymbol{\omega} \in \mathbb{R}^r : \epsilon < \|\boldsymbol{\omega}\| < 1\}$ is a finite positive measure. Lemma 3.18 asserts that the mapping $K^\epsilon : \mathbb{R}^r \rightarrow \mathbb{R}$ defined by

$$K^\epsilon(\mathbf{x}) = \int_{\epsilon < \|\boldsymbol{\omega}\| < 1} e^{i\langle \boldsymbol{\omega} | \mathbf{x} \rangle} d\boldsymbol{\omega} \quad (3.62)$$

is positive definite. It is an isotropic covariance kernel: K^ϵ only depends on \mathbf{x} through its norm $\|\mathbf{x}\|$. Let Σ_θ^ϵ be the correlation matrix corresponding with K^ϵ : its (i, i') -th element is $K^\epsilon\left(\frac{\mathbf{x}^{(i)} - \mathbf{x}^{(i')}}{\theta}\right)$. The Lebesgue measure on $\{\boldsymbol{\omega} \in \mathbb{R}^r : \epsilon < \|\boldsymbol{\omega}\| < 1\}$ is absolutely continuous with respect to the Lebesgue measure on \mathbb{R}^r , so Lemma 3.18 also asserts that Σ_θ^ϵ is positive definite.

Moreover, for every nonnegative integer k , $\int \|\boldsymbol{\omega}\|^k \mathbf{1}_{\epsilon < \|\boldsymbol{\omega}\| < 1} d\boldsymbol{\omega}$ is smaller than the mass of the Lebesgue measure on $\{\boldsymbol{\omega} \in \mathbb{R}^r : \epsilon < \|\boldsymbol{\omega}\| < 1\}$. The Maclaurin series of K^ϵ has therefore infinite radius of convergence. For any nonnegative integer k , denote by $\mathbf{D}^{(k)}$ the $n \times n$ matrix with (i, i') -th element $\left\|\mathbf{x}^{(i)} - \mathbf{x}^{(i')}\right\|^{2k}$. Because the Maclaurin series has infinite radius of convergence, Σ_θ^ϵ is equal to its asymptotic expansion regardless of the value of θ . There exist real numbers a_k ($k \in \mathbb{N}$) such that for all $\theta \in (0, +\infty)$, $\Sigma_\theta^\epsilon = \sum_{k=0}^{\infty} a_k \theta^{-2k} \mathbf{D}^{(k)}$. Because Σ_θ^ϵ is positive definite, Lemma 3.13 ensures there exists a nonnegative integer N such that the vector space $\cap_{k=0}^N \text{Ker } a_k \mathbf{D}^{(k)}$ is trivial.

Applying Lemma 3.14 then yields that when $\theta \rightarrow +\infty$ $\left\|(\Sigma_\theta^\epsilon)^{-1}\right\| = O(\theta^{2N})$. Because the greatest eigenvalue of a positive definite matrix is the smallest eigenvalue of its inverse, this implies the existence of a constant $c_\epsilon > 0$ such that when θ is large enough

$$\min_{\boldsymbol{\xi} \in \mathbb{R}^n, \|\boldsymbol{\xi}\|=1} \boldsymbol{\xi}^\top \Sigma_\theta^\epsilon \boldsymbol{\xi} \geq c_\epsilon \theta^{-2N}. \quad (3.63)$$

So when θ is large enough, for every $\boldsymbol{\xi} \in \mathbb{R}^n$,

$$\int \mathbf{1}_{\epsilon < \|\boldsymbol{\omega}\| < 1} \left| \sum_{j=1}^n \xi_j e^{i \langle \boldsymbol{\omega}, \frac{\boldsymbol{w}^{(j)}}{\theta} \rangle} \right|^2 d\boldsymbol{\omega} \geq c_\epsilon \|\boldsymbol{\xi}\|^2 \theta^{-2N}. \quad (3.64)$$

This provides a lower bound for $|J_\theta^1(\boldsymbol{\xi})|$:

$$\theta |J_\theta^1(\boldsymbol{\xi})| \geq \tilde{m}_\epsilon c_\epsilon \|\boldsymbol{\xi}\|^2 \theta^{-r-2N}. \quad (3.65)$$

Besides, we have when θ is large enough, for every $\boldsymbol{\xi} \in \mathbb{R}^n$,

$$\begin{aligned} \theta |J_\theta^2(\boldsymbol{\xi})| &\leq n^2 \|\boldsymbol{\xi}\|^2 \int_{\|\mathbf{s}\|>1} (\theta \|\mathbf{s}\|)^{\nu - \frac{r}{2} - \frac{1}{2}} \exp(-\theta \|\mathbf{s}\|) d\mathbf{s} \\ &= n^2 \|\boldsymbol{\xi}\|^2 \frac{2\pi^{\frac{r-1}{2}}}{\Gamma\left(\frac{r-1}{2}\right)} \int_1^{+\infty} (\theta t)^{\nu - \frac{r}{2} - \frac{1}{2}} \exp(-\theta t) t^{r-1} dt \\ &\leq n^2 \|\boldsymbol{\xi}\|^2 \frac{2\pi^{\frac{r-1}{2}}}{\Gamma\left(\frac{r-1}{2}\right)} \Gamma\left(\nu + \frac{r-1}{2}\right) \theta^{\nu - \frac{r}{2} - \frac{1}{2}} \exp(-(\theta - 1)). \end{aligned} \quad (3.66)$$

From Equations (3.65) and (3.66), we gather that when $\theta \rightarrow +\infty$, $\sup_{\|\boldsymbol{\xi}\|=1} |J_\theta^2(\boldsymbol{\xi})| = o(\inf_{\|\boldsymbol{\xi}\|=1} |J_\theta^1(\boldsymbol{\xi})|)$, so for any $\lambda > 1$, when θ is large enough, $-\frac{d}{d\theta} J_\theta^{RQ}(\boldsymbol{\xi}) \leq -\lambda J_\theta^1(\boldsymbol{\xi})$.

Denote by $(r - 2\nu)_+$ the quantity $\max(0, r - 2\nu)$. Then, combining Equations (3.57) and (3.59), there exists $a'_{r,\nu} > 0$ such that

$$\begin{aligned} & -\theta \frac{d}{d\theta} I_\theta^{RQ}(\boldsymbol{\xi}) \\ & \leq \lambda \int_{\|\mathbf{s}\| \leq 1} \max(a'_{r,\nu}, (\lambda + (r - 2\nu)_+) \theta \|\mathbf{s}\|) (\theta \|\mathbf{s}\|)^{\nu - \frac{r}{2}} \mathcal{K}_{\nu - \frac{r}{2}}(\theta \|\mathbf{s}\|) \left| \sum_{j=1}^n \xi_j e^{i\langle \mathbf{s} | \mathbf{x}^{(j)} \rangle} \right|^2 d\mathbf{s} \\ & \leq \lambda \max(a'_{r,\nu}, (\lambda + (r - 2\nu)_+) \theta) \int_{\|\mathbf{s}\| \leq 1} (\theta \|\mathbf{s}\|)^{\nu - \frac{r}{2}} \mathcal{K}_{\nu - \frac{r}{2}}(\theta \|\mathbf{s}\|) \left| \sum_{j=1}^n \xi_j e^{i\langle \mathbf{s} | \mathbf{x}^{(j)} \rangle} \right|^2 d\mathbf{s}. \end{aligned} \quad (3.67)$$

When θ is large enough, $a'_{r,\nu} \leq \lambda \theta$ so

$$-\theta \frac{d}{d\theta} I_\theta^{RQ}(\boldsymbol{\xi}) \leq \lambda (\lambda + (r - 2\nu)_+) \theta I_\theta^{RQ}(\boldsymbol{\xi}). \quad (3.68)$$

From this, we obtain that for any non-null vector $\boldsymbol{\xi} \in \mathbb{R}^n$, for any $\lambda > 1$, provided θ is large enough,

$$0 < -\frac{d}{d\theta} I_\theta^{RQ}(\boldsymbol{\xi}) \leq \lambda (\lambda + (r - 2\nu)_+) I_\theta^{RQ}(\boldsymbol{\xi}). \quad (3.69)$$

Combining Equations (3.48), (3.53) and (3.69) yields the result. \square

Lemma 3.22. *For the Squared Exponential kernel, $\mathbf{F}_\theta := r\theta^{-1}\Sigma_\theta - \frac{d}{d\theta}\Sigma_\theta$ is a symmetric positive definite matrix. If θ is large enough, it verifies $\forall \boldsymbol{\xi} \in \mathbb{R}^n$, $\boldsymbol{\xi}^\top \mathbf{F}_\theta \boldsymbol{\xi} \leq \theta \boldsymbol{\xi}^\top \Sigma_\theta \boldsymbol{\xi}$.*

Proof. For any $\theta \in (0, +\infty)$ and any $\boldsymbol{\xi} \in \mathbb{R}^n$,

$$\frac{d}{d\theta} I_\theta^{SE}(\boldsymbol{\xi}) = \int_{\mathbb{R}^r} -\frac{\theta \|\mathbf{s}\|^2}{2} \exp\left(-\frac{\theta^2 \|\mathbf{s}\|^2}{4}\right) \left| \sum_{j=1}^n \xi_j e^{i\langle \mathbf{s} | \mathbf{x}^{(j)} \rangle} \right|^2 d\mathbf{s} \quad (3.70)$$

So \mathbf{F}_θ is positive definite.

Similarly to the Rational Quadratic case (cf. proof of Lemma 3.21), one can show that for any $\lambda > 1$, for large enough θ and for any non-null vector $\boldsymbol{\xi} \in \mathbb{R}^n$,

$$0 < -\frac{d}{d\theta} I_\theta^{SE}(\boldsymbol{\xi}) \leq \lambda \int_{\|\mathbf{s}\| \leq 1} \frac{\theta \|\mathbf{s}\|^2}{2} \exp\left(-\frac{\theta^2 \|\mathbf{s}\|^2}{4}\right) \left| \sum_{j=1}^n \xi_j e^{i\langle \mathbf{s} | \mathbf{x}^{(j)} \rangle} \right|^2 d\mathbf{s} \leq \frac{\lambda}{2} \theta I_\theta^{SE}(\boldsymbol{\xi}). \quad (3.71)$$

Combining Equations (3.48), (3.53) and (3.71) yields the result. \square

Appendix 3.D Asymptotic study of the correlation matrix Σ_θ

Rational Quadratic and Squared Exponential kernels

For all $\nu > 0$, the series expansion of the mapping $x \mapsto (1 + x)^{-\nu}$ at $x = 0$ has radius of convergence 1. Moreover, the series expansion of the exponential function has infinite radius

of convergence. From these facts follows that when θ is large enough, if a Rational Quadratic kernel or a Squared Exponential kernel is used,

$$\Sigma_\theta = \sum_{k=0}^{\infty} \frac{a_k}{\theta^{2k}} \mathbf{D}^{(k)}. \quad (3.72)$$

In the above expression, for every k , $\mathbf{D}^{(k)}$ is the $n \times n$ matrix with (i, i') -th element $\|\mathbf{x}^{(i)} - \mathbf{x}^{(i')}\|^{2k}$ and a_k is a non-null real number. To be precise, $a_k = (-1)^k \left(\prod_{l=0}^k (\nu + l) \right) / k!$ for Rational Quadratic kernels and $a_k = (-1)^k / k!$ for the Squared Exponential kernel.

Equation (3.72) implies

$$\mathbf{W}^\top \Sigma_\theta \mathbf{W} = \sum_{k=0}^{\infty} \frac{a_k}{\theta^{2k}} \mathbf{W}^\top \mathbf{D}^{(k)} \mathbf{W}. \quad (3.73)$$

Σ_θ is positive definite and the kernel of \mathbf{W} is trivial so $\mathbf{W}^\top \Sigma_\theta \mathbf{W}$ is positive definite. Let k_1 be the smallest nonnegative integer such that $\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W}$ is non-null. Define $\mathbf{D} := a_{k_1} \mathbf{D}^{(k_1)}$. If $\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W}$ is nonsingular, then define $k_2 := k_1 + 1$ and $\mathbf{D}^* := a_{k_2} \mathbf{D}^{(k_2)}$. If $\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W}$ is singular, then because $\mathbf{W}^\top \Sigma_\theta \mathbf{W}$ is nonsingular, there must exist an integer $k > k_1$ such that $\mathbf{W}^\top \mathbf{D}^{(k)} \mathbf{W}$ is non-null. Then let k_2 be the smallest of them and define $\mathbf{D}^* := a_{k_2} \mathbf{D}^{(k_2)}$. Now, define the mappings $g(\theta) = \theta^{-2k_1}$ and $g^*(\theta) = \theta^{-2(k_2 - k_1)}$.

Finally, define

$$\mathbf{R}_g(\theta) = g(\theta)^{-1} \sum_{k=k_2+1}^{\infty} \frac{a_k}{\theta^{2k}} \mathbf{W}^\top \mathbf{D}^{(k)} \mathbf{W}. \quad (3.74)$$

Notice that $\|\mathbf{R}_g(\theta)\| = o(g^*(\theta))$ and that $\|\frac{d}{d\theta} \mathbf{R}_g(\theta)\| = o(g^{*'}(\theta))$.

Matérn kernels with noninteger smoothness ν

If a Matérn kernel with noninteger smoothness $\nu > 0$ (whether greater or smaller than 1) is used, we can write Σ_θ as [Abramowitz and Stegun, 1964] (9.6.2. and 9.6.10.):

$$\Sigma_\theta = \sum_{k=0}^{\lfloor \nu \rfloor} \frac{a_k}{\theta^{2k}} \mathbf{D}^{(k)} + \frac{a_\nu}{\theta^{2\nu}} \mathbf{D}^{(\nu)} + \mathbf{R}(\theta). \quad (3.75)$$

Like in the case of Rational Quadratic and Squared Exponential kernels, for every k , $\mathbf{D}^{(k)}$ is the $n \times n$ matrix with (i, i') -th element $\|\mathbf{x}^{(i)} - \mathbf{x}^{(i')}\|^{2k}$. The a_k 's, of course, are different: $a_k = (-1)^k \Gamma(\nu - k) \nu^k / (k! \Gamma(\nu))$. Moreover, $\mathbf{D}^{(\nu)}$ is the $n \times n$ matrix with (i, i') -th element $\|\mathbf{x}^{(i)} - \mathbf{x}^{(i')}\|^{2\nu}$, $a_\nu = \Gamma(-\nu) \nu^\nu / \Gamma(\nu)$ and \mathbf{R} is a differentiable mapping from $(0, +\infty)$ to the set of real $n \times n$ matrices \mathcal{M}_n such that $\|\mathbf{R}(\theta)\| = O(\theta^{-2(\lfloor \nu \rfloor + 1)})$ and $\|\frac{d}{d\theta} \mathbf{R}(\theta)\| = O(\theta^{-2(\lfloor \nu \rfloor + 1) - 1})$. Lemma (3.19) implies that when θ is large enough $\Sigma_\theta - \mathbf{R}(\theta)$ is positive definite.

Equation (3.75) implies

$$\mathbf{W}^\top \Sigma_\theta \mathbf{W} = \sum_{k=0}^{\lfloor \nu \rfloor} \frac{a_k}{\theta^{2k}} \mathbf{W}^\top \mathbf{D}^{(k)} \mathbf{W} + \frac{a_\nu}{\theta^{2\nu}} \mathbf{W}^\top \mathbf{D}^{(\nu)} \mathbf{W} + \mathbf{W}^\top \mathbf{R}(\theta) \mathbf{W}. \quad (3.76)$$

When θ is large enough, $\Sigma_\theta - \mathbf{R}(\theta)$ is positive definite. Since the kernel of \mathbf{W} is trivial, when θ is large enough, $\mathbf{W}^\top \Sigma_\theta \mathbf{W} - \mathbf{W}^\top \mathbf{R}(\theta) \mathbf{W}$ is positive definite. If it exists, let k_1 be

the smallest nonnegative integer smaller than ν such that $\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W}$ is non-null and define $\mathbf{D} := a_{k_1} \mathbf{D}^{(k_1)}$ and $g(\theta) := \theta^{-2k_1}$. If not, then define $\mathbf{D} := a_\nu \mathbf{D}^{(\nu)}$ and $g(\theta) := \theta^{-2\nu}$. Either way $\mathbf{W}^\top \mathbf{D} \mathbf{W}$ is non-null.

If k_1 exists and $\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W}$ is nonsingular, then define $k_2 := k_1 + 1$ if $k_1 < \lfloor \nu \rfloor$ and $k_2 = \nu$ if $k_1 = \lfloor \nu \rfloor$. Then define $\mathbf{D}^* := a_{k_2} \mathbf{D}^{(k_2)}$ and $g^*(\theta) := g(\theta)^{-1} \theta^{-2k_2}$.

If k_1 exists and $\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W}$ is singular, then there must exist $k \in \llbracket k_1 + 1, \lfloor \nu \rfloor \rrbracket \cup \{\nu\}$ such that $\mathbf{W}^\top \mathbf{D}^{(k)} \mathbf{W}$ is non-null. Let k_2 be the smallest number among all such k . Define $\mathbf{D}^* := a_{k_2} \mathbf{D}^{(k_2)}$ and $g^*(\theta) := g(\theta)^{-1} \theta^{-2k_2}$.

If k_1 does not exist, then $\mathbf{W}^\top \mathbf{D}^{(\nu)} \mathbf{W}$ is necessarily nonsingular. Define \mathbf{D}^* as the null $n \times n$ matrix and $g^*(\theta) = g(\theta)^{-1} \theta^{-2\nu - (\lfloor \nu \rfloor + 1 - \nu)}$.

Finally, define

$$\mathbf{R}_g(\theta) := g(\theta)^{-1} \left(\mathbf{W}^\top \Sigma_\theta \mathbf{W} - g(\theta) \mathbf{W}^\top \mathbf{D} \mathbf{W} - g(\theta) g^*(\theta) \mathbf{W}^\top \mathbf{D}^* \mathbf{W} \right). \quad (3.77)$$

In all situations, $\|\mathbf{R}_g(\theta)\| = o(g^*(\theta))$ and $\|\frac{d}{d\theta} \mathbf{R}_g(\theta)\| = o(g^*(\theta))$.

Matérn kernels with integer smoothness ν

Finally, if a Matérn kernel with integer smoothness ν is used, we can write Σ_θ as [Abramowitz and Stegun, 1964] (9.6.11.):

$$\Sigma_\theta = \sum_{k=0}^{\nu-1} \frac{a_k}{\theta^{2k}} \mathbf{D}^{(k)} + \tilde{a}_\nu \left(\frac{\log(\theta)}{\theta^{2\nu}} \mathbf{D}^{(\nu)} + \frac{1}{\theta^{2\nu}} \tilde{\mathbf{D}}^{(\nu)} \right) + \mathbf{R}(\theta). \quad (3.78)$$

a_k and $\mathbf{D}^{(k)}$ ($k \in \llbracket 0, \nu - 1 \rrbracket$) and $\mathbf{D}^{(\nu)}$ have the same definitions as for Matérn kernels with noninteger smoothness $\nu > 1$.

$\tilde{a}_\nu = (-1)^\nu 2\nu^\nu / (\nu - 1)!$ and $\tilde{\mathbf{D}}^{(\nu)}$ is the $n \times n$ matrix with null diagonal and (i, i') -th element ($i \neq i'$) given by $\|\mathbf{x}^{(i)} - \mathbf{x}^{(i')}\|^{2\nu} \left\{ -0.5 \log(\|\mathbf{x}^{(i)} - \mathbf{x}^{(i')}\|^2) - 0.5 \log(\nu) - 2\gamma + \sum_{l=1}^{\nu} l^{-1} \right\}$, where γ is Euler's constant.

Finally, \mathbf{R} is a differentiable mapping from $(0, +\infty)$ to the set of real $n \times n$ matrices \mathcal{M}_n such that $\|\mathbf{R}(\theta)\| = O(\log(\theta)\theta^{-2(\nu+1)})$ and $\|\frac{d}{d\theta} \mathbf{R}(\theta)\| = O(\log(\theta)\theta^{-2(\nu+1)-1})$.

Equation (3.78) implies

$$\mathbf{W}^\top \Sigma_\theta \mathbf{W} = \sum_{k=0}^{\nu-1} \frac{a_k}{\theta^{2k}} \mathbf{W}^\top \mathbf{D}^{(k)} \mathbf{W} + \frac{\log(\theta)}{\theta^{2\nu}} \tilde{a}_\nu \mathbf{W}^\top \mathbf{D}^{(\nu)} \mathbf{W} + \frac{\tilde{a}_\nu}{\theta^{2\nu}} \mathbf{W}^\top \tilde{\mathbf{D}}^{(\nu)} \mathbf{W} + \mathbf{W}^\top \mathbf{R}(\theta) \mathbf{W}. \quad (3.79)$$

When θ is large enough, $\Sigma_\theta - \mathbf{R}(\theta)$ is positive definite. Since the kernel of \mathbf{W} is trivial, when θ is large enough, $\mathbf{W}^\top \Sigma_\theta \mathbf{W} - \mathbf{W}^\top \mathbf{R}(\theta) \mathbf{W}$ is positive definite. If it exists, let k_1 be the smallest nonnegative integer smaller or equal to ν such that $\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W}$ is non-null and define $\mathbf{D} := a_{k_1} \mathbf{D}^{(k_1)}$ and $g(\theta) = \theta^{-2k_1}$ ($k_1 < \nu$) or $\mathbf{D} := \tilde{a}_\nu \mathbf{D}^{(\nu)}$ and $g(\theta) := \log(\theta)\theta^{-2\nu}$ ($k_1 = \nu$). If not, then define $\mathbf{D} := \tilde{a}_\nu \tilde{\mathbf{D}}^{(\nu)}$ and $g(\theta) := \theta^{-2\nu}$. Either way $\mathbf{W}^\top \mathbf{D} \mathbf{W}$ is non-null.

If $\mathbf{W}^\top \mathbf{D} \mathbf{W}$ is nonsingular, then

- either k_1 exists and $k_1 < \nu - 1$, in which case define $\mathbf{D}^* := a_{k_1+1} \mathbf{D}^{(k_1+1)}$ and $g^*(\theta) := \theta^{-2}$;

- or k_1 exists and $k_1 = \nu - 1$, in which case define $\mathbf{D}^* := \tilde{a}_\nu \mathbf{D}^{(\nu)}$ and $g^*(\theta) := \log(\theta)\theta^{-2}$;
- or k_1 exists and $k_1 = \nu$, in which case define $\mathbf{D}^* := \tilde{a}_\nu \tilde{\mathbf{D}}^\nu$ and $g^*(\theta) := \log(\theta)^{-1}$;
- or k_1 does not exist, in which case define \mathbf{D}^* as the null $n \times n$ matrix and $g^*(\theta) := \theta^{-1}$.

If $\mathbf{W}^\top \mathbf{D} \mathbf{W}$ is singular, then k_1 necessarily exists:

- either $k_1 < \nu$. Then there are two possibilities. The first is that there exists a smallest integer $k_2 \in \llbracket k_1 + 1, \nu \rrbracket$ such that $\mathbf{W}^\top \mathbf{D}^{(k_2)} \mathbf{W}$ is non-null, in which case define $\mathbf{D}^* := a_{k_2} \mathbf{D}^{(k_2)}$ and $g^*(\theta) := \theta^{-2(k_2 - k_1)}$ ($k_2 < \nu$) or $\mathbf{D}^* := \tilde{a}_\nu \mathbf{D}^{(\nu)}$ and $g^*(\theta) := \log(\theta)\theta^{-2(\nu - k_1)}$ ($k_2 = \nu$). The second is that no such k_2 exists, but then $\mathbf{W}^\top \tilde{\mathbf{D}}^{(\nu)} \mathbf{W}$ is necessarily non-null, so define $\mathbf{D}^* := \tilde{a}_\nu \tilde{\mathbf{D}}^{(\nu)}$ and $g^*(\theta) := \theta^{-2(\nu - k_1)}$.
- or $k_1 = \nu$. Then $\mathbf{W}^\top \tilde{\mathbf{D}}^{(\nu)} \mathbf{W}$ is necessarily non-null, so define $\mathbf{D}^* := \tilde{a}_\nu \tilde{\mathbf{D}}^{(\nu)}$ and $g^*(\theta) := \log(\theta)^{-1}$.

Finally, define

$$\mathbf{R}_g(\theta) := g(\theta)^{-1} \left(\mathbf{W}^\top \Sigma_\theta \mathbf{W} - g(\theta) \mathbf{W}^\top \mathbf{D} \mathbf{W} - g(\theta) g^*(\theta) \mathbf{W}^\top \mathbf{D}^* \mathbf{W} \right). \quad (3.80)$$

In all situations, $\|\mathbf{R}_g(\theta)\| = o(g^*(\theta))$ and $\|\frac{d}{d\theta} \mathbf{R}_g(\theta)\| = o(g^{*'}(\theta))$.

Proof of Lemma 3.8

For Rational Quadratic and Squared Exponential kernels, Equation (3.72) implies thanks to Lemma 3.13 that there exists $k' \in \mathbb{N}$ such that $\cap_{k=0}^{k'} \text{Ker} \left(\mathbf{W}^\top \mathbf{D}^{(k)} \mathbf{W} \right)$ is the trivial vector space and $\cap_{0 \leq k < k'} \text{Ker} \left(\mathbf{W}^\top \mathbf{D}^{(k)} \mathbf{W} \right)$ is a non-trivial vector space (if $k' = 0$, the intersection is done over an empty index set, so we take it to be \mathbb{R}^{n-p} by convention).

Lemma 3.14 implies that there exists a constant $c_{k'} > 0$ such that for large enough θ , $v_{n-p}(\theta) \geq c_{k'} \theta^{-2k'}$. Thanks to Lemma 3.16, there exists a hyperplane \mathcal{H}_{n-p} of \mathbb{R}^{n-p} such that for every $\mathbf{y}' \in \mathbb{R}^{n-p} \setminus \mathcal{H}_{n-p}$, there exists $c_{\mathbf{y}'} > 0$ such that for large enough θ ,

$$(\mathbf{y}')^\top \left(\mathbf{W}^\top \Sigma_\theta \mathbf{W} \right)^{-1} \mathbf{y}' \geq c_{\mathbf{y}'} \left\| \left(\mathbf{W}^\top \Sigma_\theta \mathbf{W} \right)^{-1} \right\|. \quad (3.81)$$

So for every $\mathbf{y} \in \mathbb{R}^n$ such that $\mathbf{W}^\top \mathbf{y} \in \mathbb{R}^{n-p} \setminus \mathcal{H}_{n-p}$, there exists $c_{\mathbf{y}} > 0$ such that for large enough θ

$$\mathbf{y}^\top \mathbf{W} \left(\mathbf{W}^\top \Sigma_\theta \mathbf{W} \right)^{-1} \mathbf{W}^\top \mathbf{y} \geq c_{\mathbf{y}} \left\| \left(\mathbf{W}^\top \Sigma_\theta \mathbf{W} \right)^{-1} \right\|. \quad (3.82)$$

Because the matrix \mathbf{W}^\top has full row rank, the vector space of all $\mathbf{v} \in \mathbb{R}^n$ such that $\mathbf{W}^\top \mathbf{v} \in \mathcal{H}_{n-p}$ is included within a hyperplane \mathcal{H}_n of \mathbb{R}^n , so for every $\mathbf{y} \in \mathbb{R}^n \setminus \mathcal{H}_n$, there exists $c_{\mathbf{y}} > 0$ such that for large θ the above equation holds.

For Matérn kernels with noninteger smoothness $\nu > 0$ (resp. with integer smoothness $\nu > 0$), Equation (3.75) (resp. Equation (3.78)) allows a similar argument. Indeed, Lemma 3.19 asserts that $\|\Sigma_\theta^{-1}\| = O(\theta^{2\nu})$, so $\cap_{k=0}^{\lfloor \nu \rfloor} \text{Ker} \left(\mathbf{W}^\top \mathbf{D}^{(k)} \mathbf{W} \right) \cap \text{Ker} \left(\mathbf{W}^\top \mathbf{D}^{(\nu)} \mathbf{W} \right)$ (resp. $\cap_{k=0}^{\nu} \text{Ker} \left(\mathbf{W}^\top \mathbf{D}^{(k)} \mathbf{W} \right) \cap \text{Ker} \left(\mathbf{W}^\top \tilde{\mathbf{D}}^{(\nu)} \mathbf{W} \right)$) is necessarily the trivial vector space.

Appendix 3.E Details of the proof of Theorem 3.9

Here the last part of the proof of Theorem 3.9 is given in detail.

Rational Quadratic, Squared Exponential and Matérn ($\nu \in [1, +\infty) \setminus \mathbb{Z}_+$) kernels

Let us first tackle the case of Rational Quadratic and Squared Exponential kernels and of Matérn kernels with noninteger smoothness $\nu > 1$.

In case 1. (a), Lemma 3.14 yields $\tilde{w}(\theta) = O(g^{*'}(\theta))$.

This implies $\pi(\theta) = O(g^{*'}(\theta)) = O(\theta^{-2l-1})$, so the reference prior is proper. Given the likelihood function is bounded (cf. Equation 3.12), the reference posterior is proper as well.

In case 1. (b), Lemma 3.14 yields $\tilde{w}(\theta) = O(g^{*'}(\theta)g^*(\theta)^{-1})$.

This implies $\pi(\theta) = O(g^{*'}(\theta)g^*(\theta)^{-1}) = O(\theta^{-1})$. Moreover, $v_{n-p}(\theta) = O(g(\theta)g^*(\theta))$. As the rank of $\mathbf{W}^\top \mathbf{D} \mathbf{W}$ is at least one, $v_1(\theta)^{-1} = O(g(\theta)^{-1})$. Gathering all this, $v_{n-p}(\theta)/v_1(\theta) = O(g^*(\theta))$, so Equation (3.12) implies $L(\mathbf{y}|\theta) = O(g^*(\theta))^{1/2} = O(\theta^{-l})$. The reference posterior is then proportional to $L(\mathbf{y}|\theta)\pi(\theta) = O(\theta^{-l-1})$ and is proper.

In case 2., we must distinguish between Matérn kernels and the others. For Matérn kernels with noninteger smoothness $\nu > 1$, Proposition 3.7 asserts that the reference prior is $O(\theta^{-1})$ so the argument used in case 1.(b) still holds. For Rational Quadratic and Squared Exponential kernels, Equation (3.72) implies

$$\mathbf{W}^\top \Sigma_\theta \mathbf{W} = \sum_{k=0}^{\infty} \frac{a_k}{\theta^{2k}} \mathbf{W}^\top \mathbf{D}^{(k)} \mathbf{W}. \quad (3.83)$$

Let k_1 be the smallest nonnegative integer such that $\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W}$ is not the null matrix. Then

$$g(\theta) \mathbf{W}^\top \mathbf{D} \mathbf{W} = a_{k_1} \theta^{-2k_1} \mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W}. \quad (3.84)$$

and for some integer $k_2 > k_1$,

$$g(\theta)g^*(\theta) \mathbf{W}^\top \mathbf{D}^* \mathbf{W} = a_{k_2} \theta^{-2k_2} \mathbf{W}^\top \mathbf{D}^{(k_2)} \mathbf{W}. \quad (3.85)$$

Things are easiest if $\mathbf{W}^\top \mathbf{D}^{(k_1+1)} \mathbf{W}$ is null, because then $k_2 > k_1 + 1$. Since we are dealing with case 2., $\text{Ker}(\mathbf{W}^\top \mathbf{D}^* \mathbf{W}) \cap \text{Ker}(\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W})$ is not the trivial vector space so Equation (3.83) yields $v_{n-p}(\theta) = O(\theta^{-2(k_2+1)}) = O(\theta^{-2(k_1+3)})$. Besides, the smallest eigenvalue of $(\mathbf{W}^\top \Sigma_\theta \mathbf{W})^{-1}$ verifies $v_1(\theta)^{-1} = O(\theta^{-2k_1})$ so $v_{n-p}(\theta)/v_1(\theta) = O(\theta^{-6})$. Recall Proposition 3.7 asserts that $\pi(\theta) = O(\theta)$. The reference posterior is proportional to $L(\mathbf{y}|\theta)\pi(\theta) = O(\theta^{-2})$ and thus proper.

In the following, assume $\mathbf{W}^\top \mathbf{D}^{(k_1+1)} \mathbf{W}$ is not null. Then $k_2 = k_1 + 1$.

If we assume that $\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W}$ has rank greater or equal to 2, then a similar reasoning can be applied. The two smallest eigenvalues of $(\mathbf{W}^\top \Sigma_\theta \mathbf{W})^{-1}$ verify $v_1(\theta)^{-1} = O(\theta^{-2k_1})$ and $v_2(\theta)^{-1} = O(\theta^{-2k_1})$. $v_{n-p}(\theta) = O(\theta^{-2(k_2+1)}) = O(\theta^{-2(k_1+2)})$. From this we obtain $v_{n-p}(\theta)/v_1(\theta) = O(\theta^{-4})$ and $v_{n-p}(\theta)/v_2(\theta) = O(\theta^{-4})$. Equation (3.12) then implies $L(\mathbf{y}|\theta) = O(\theta^{-2-2})$. The reference posterior is proportional to $L(\mathbf{y}|\theta)\pi(\theta) = O(\theta^{-3})$ and thus proper.

Now, assume that $\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W}$ has rank 1. We need to distinguish between two possibilities: either $\text{Ker}(\mathbf{W}^\top \mathbf{D}^{(k_1+2)} \mathbf{W}) \cap \text{Ker}(\mathbf{W}^\top \mathbf{D}^{(k_1+1)} \mathbf{W}) \cap \text{Ker}(\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W})$ is the trivial vector space, or it is not. If it is not, Equation (3.83) yields $v_{n-p}(\theta) = O(\theta^{-2(k_1+3)})$ and the conclusion is the same as in the case where $\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W}$ is null.

Let us now deal with the situation where $\text{Ker}(\mathbf{W}^\top \mathbf{D}^{(k_1+2)} \mathbf{W}) \cap \text{Ker}(\mathbf{W}^\top \mathbf{D}^{(k_1+1)} \mathbf{W}) \cap \text{Ker}(\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W})$ is the trivial vector space. Two further subcases must be distinguished here: either $\text{Ker}(\mathbf{W}^\top \mathbf{D}^{(k_1+1)} \mathbf{W}) \cap \text{Ker}(\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W})$ is equal to $\text{Ker}(\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W})$, or it is strictly included within $\text{Ker}(\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W})$.

If it is strictly included, then Lemma 3.17 is applicable and $v_1(\theta)^{-1} = O(\theta^{-2k_1})$ and $v_2(\theta)^{-1} = O(\theta^{-2(k_1+1)})$. Because this is case 2., $\text{Ker}(\mathbf{W}^\top \mathbf{D}^{(k_1+1)} \mathbf{W}) \cap \text{Ker}(\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W})$ is not the trivial vector space, so Equation (3.83) yields $v_{n-p}(\theta) = O(\theta^{-2(k_1+2)})$. From there we obtain $v_{n-p}(\theta)/v_1(\theta) = O(\theta^{-4})$ and $v_{n-p}(\theta)/v_2(\theta) = O(\theta^{-2})$. Equation (3.12) then implies $L(\mathbf{y}|\theta) = O(\theta^{-3})$, so the reference posterior is proper.

In the second subcase, $\text{Ker}(\mathbf{W}^\top \mathbf{D}^{(k_1+1)} \mathbf{W}) \cap \text{Ker}(\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W}) = \text{Ker}(\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W})$. Since $\text{Ker}(\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W})$ is a hyperplane of \mathbb{R}^{n-p} , this implies that $\text{Ker}(\mathbf{W}^\top \mathbf{D}^{(k_1+1)} \mathbf{W})$ is the same hyperplane: $\mathbf{W}^\top \mathbf{D}^{(k_1+1)} \mathbf{W}$ and $\mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W}$ are symmetric matrices of rank 1 with the same kernel. So there exists $b \in \mathbb{R} \setminus \{0\}$ such that $\mathbf{W}^\top \mathbf{D}^{(k_1+1)} \mathbf{W} = b \mathbf{W}^\top \mathbf{D}^{(k_1)} \mathbf{W}$. If we redefine $g(\theta) := a_{k_1} \theta^{-2k_1} + a_{k_1+1} b \theta^{-2(k_1+1)}$, $g^*(\theta) := a_{k_1+2} \theta^{-2(k_1+2)} g(\theta)^{-1}$, $\mathbf{D}^* := \mathbf{D}^{(k_1+2)}$ and $\mathbf{R}_g(\theta) := g(\theta)^{-1} \sum_{k=k_1+3}^{\infty} a_k \theta^{-2k}$, the situation is similar to case 1.(b), except that $g^*(\theta)$ is not necessarily of the form θ^{-2l} .

However, $g^{*\prime}(\theta) = g^*(\theta)(-2(k_1+2)\theta^{-1} - g'(\theta)g(\theta)^{-1})$ and $g'(\theta)g(\theta)^{-1} = O(\theta^{-1})$, which implies $g^{*\prime}(\theta)g^*(\theta)^{-1} = O(\theta^{-1})$. In addition, $g^*(\theta) = O(\theta^{-2})$. Therefore the arguments of the study of case 1.(b) apply: $\pi(\theta) = O(\theta^{-1})$ and $L(\mathbf{y}|\theta) = O(g^*(\theta)^{1/2}) = O(\theta^{-1})$. The reference posterior is proportional to $L(\mathbf{y}|\theta)\pi(\theta) = O(\theta^{-2})$: it is proper. This particular subcase, because it is analogous to case 1.(b) is called ‘‘special’’. All other subcases of case 2. collectively form the ‘‘usual’’ case.

Matérn ($\nu \in \mathbb{Z}_+$) kernels

We now address the case where the correlation kernel is Matérn with integer smoothness ν . The proof strategy remains the same as for the other kernels, but the execution is a little trickier.

It still relies on the asymptotic expansion of Σ_θ . For Matérn kernels with integer smoothness ν , the decomposition is detailed in Appendix 3.D.

First, assume either \mathbf{D} is *not* proportional to $\mathbf{D}^{(\nu)}$ or \mathbf{D}^* is *not* proportional to $\tilde{\mathbf{D}}^{(\nu)}$. In Equation (3.13), $g^*(\theta)$ may be $\theta^{-2l} \log(\theta)$ instead of θ^{-2l} . Then its derivative is $g^{*\prime}(\theta) = \theta^{-2l-1}(1 - 2l \log(\theta))$. In case 1.(a), the reference prior (and posterior) is $O(g^{*\prime}(\theta)) = O(\theta^{-2l-1} \log(\theta))$ and thus proper. It is useless to distinguish cases 1.(b) and 2. because thanks to Proposition 3.7, the reference prior is $O(\theta^{-1})$. In either case, the rank of $\mathbf{W}^\top \mathbf{D} \mathbf{W}$ is at least one, so $v_{n-p}(\theta)/v_1(\theta) = O(g^*(\theta))$. Equation (3.12) implies $L(\mathbf{y}|\theta) = O(g^*(\theta)^{1/2}) = O(\theta^{-l} \log(\theta)^{1/2})$, so the reference posterior is proportional to $L(\mathbf{y}|\theta)\pi(\theta) = O(\theta^{-l-1} \log(\theta)^{1/2})$ and thus proper.

Now, assume \mathbf{D} is proportional to $\mathbf{D}^{(\nu)}$ and \mathbf{D}^* is proportional to $\tilde{\mathbf{D}}^{(\nu)}$. In Equation (3.13), $g^*(\theta) = \log(\theta)^{-1}$. Its derivative is $g^{*\prime}(\theta) = -\theta^{-1} \log(\theta)^{-2}$. In case 1.(a), the reference prior (resp. posterior) is $O(g^{*\prime}(\theta)) = O(\theta^{-1} \log(\theta)^{-2})$ and is thus proper. In case 1.(b), the reference prior is $O(g^{*\prime}(\theta)g^*(\theta)^{-1}) = O(\theta^{-1} \log(\theta)^{-1})$. Besides, as the rank of $\mathbf{W}^\top \mathbf{D} \mathbf{W}$ is at

least one, Lemma 3.17 yields $v_{n-p}(\theta)/v_1(\theta) = O(g^*(\theta))$, so Equation (3.12) implies $L(\mathbf{y}|\theta) = O(g^*(\theta)^{1/2}) = O(\log(\theta)^{-1/2})$. The reference posterior is then proportional to $L(\mathbf{y}|\theta)\pi(\theta) = O(\theta^{-1} \log(\theta)^{-3/2})$: it is proper. Case 2. cannot occur because Lemma 3.19 asserts that $\|\Sigma_\theta^{-1}\| = O(\theta^{2\nu})$.

Part II

Compromise

Chapter 4

Optimal compromise between incompatible conditional probability distributions

This chapter covers Section 2 of the article Muré [2017].

Abstract

Models are often defined through conditional rather than joint distributions, but it can be difficult to check whether the conditional distributions are compatible. When they are not, we give meaning to the intuition that the stationary probability distribution of the associated Pseudo-Gibbs sampler is the optimal compromise between the conditional distributions.

Résumé

Les modèles statistiques sont souvent définis par lois conditionnelles plutôt que jointes. Il peut cependant être difficile de contrôler la compatibilité des lois conditionnelles. Dans le cas où elles ne seraient pas compatibles, nous donnons sens à l'intuition selon laquelle la distribution stationnaire du pseudo-échantillonneur de Gibbs correspondant est le compromis optimal entre les lois conditionnelles.

4.1 Introduction

Generally speaking, there are two ways to create statistical models for multiple random variables. One can either consider them simultaneously and directly define their joint distribution, or one can define a system of conditional distributions. The first approach is conceptually easier and often (but not always) leads to models with well understood properties. The second one allows for more flexibility in modeling but makes theoretical analysis more difficult.

The main problem with the second approach is that conditional distributions may not be compatible. In this context, it means that there exists no joint distribution from which the conditional distributions can all be derived. Other definitions of compatibility exist in the literature. For example, in the context of a model with a given prior distribution, Dawid

and Lauritzen [2001] examine the problem of eliciting a compatible prior distribution for a submodel. In the domain of Bayesian Networks, a probability distribution can be compatible or not with a given Directed Acyclic Graph (DAG) [Roverato, 2004]. Moreover, an abstraction (simplification) of a DAG can be compatible or not [Yet and Marsh, 2014]. In this dissertation however, the notion of compatibility concerns families of conditional distributions. A family of conditional distributions is called compatible if there exists a joint distribution that agrees with them all.

Accordingly, the problem of efficiently determining whether a given system of conditionals is compatible has received considerable attention over the years. Kuo et al. [2017], after listing previous attempts, provide probably the best solution to date. Their idea relies on the Structural Ratio Matrix which contains ratios between conditional distributions. Obviously, this solution, like its predecessors deals only with discrete conditional probability distributions.

However, even if a system contains incompatible conditional probability distributions, it does not follow that it is useless. Since Heckerman et al. [2000], practitioners have been using systems of conditional probability distributions without reference to compatibility. Gelman and Raghunathan [2001] state that “in general, reasonable-seeming conditional models will not be compatible with any single joint distribution”. It is, indeed, always possible to fire up Gibbs samplers to deal with a system of conditional distributions. Some authors use the colorful acronym PIGS for “Potentially Incompatible Gibbs Sampler” to describe such a procedure. When the conditionals are definitely known to be incompatible, the most widely used term seems to be “Pseudo-Gibbs Sampler” (PGS).

Behind the practice of PIGS is the intuition that the Gibbs sampler should converge to the joint distribution that best represents the system of conditionals. Kuo and Wang [2017] provide a detailed analysis and geometrical interpretation of the behavior of Pseudo-Gibbs Samplers for discrete conditional distributions. In particular, they show how the scanning order determines its stationary distribution. In Section 4.2 of the present chapter, we provide some theoretical foundation for the intuition that the stationary distribution of a PGS with random scanning order is, in case of existence and uniqueness, the best “compromise” between incompatible conditionals.

4.2 Optimal compromise: a general theory

Definitions and notations

In this section we introduce the concepts necessary to define the optimal compromise between potentially incompatible conditional distributions.

First, note that in this context, “conditional distribution” is really an informal way of referring to a Markov kernel.

Definition 4.1. *Let (A, \mathcal{A}) and (B, \mathcal{B}) be measurable sets. A mapping $\pi : A \times \mathcal{B} \rightarrow [0, 1]$ is called a Markov kernel if:*

1. *for all $x \in A$, $\pi(x, \cdot) : \mathcal{B} \rightarrow [0, 1]$ is a probability distribution and*
2. *for all $S \in \mathcal{B}$, $\pi(\cdot, S) : A \rightarrow [0, 1]$ is \mathcal{A} -measurable.*

We use the following notation: for every $(x, S) \in A \times \mathcal{B}$, $\pi(S|x) := \pi(x, S)$.

Let r be a positive integer and let $(\Omega_1, \mathcal{A}_1), \dots, (\Omega_r, \mathcal{A}_r)$ be measurable sets. Define $\Omega = \times_{i=1}^r \Omega_i = \Omega_1 \times \dots \times \Omega_r$ and $\mathcal{A} := \otimes_{i=1}^r \mathcal{A}_i = \mathcal{A}_1 \otimes \dots \otimes \mathcal{A}_r$.

For every $i \in \llbracket 1, r \rrbracket$, let π_i be a Markov kernel $(\times_{j \neq i} \Omega_j) \times \mathcal{A}_i \rightarrow [0, 1]$.

Intuitively (we formalize this below), every π_i should be assembled with a distribution m^i on $\otimes_{j \neq i} \mathcal{A}_j$ to create a ‘‘joint’’ distribution, that is a probability distribution on \mathcal{A} . We refer to every m^i ($i \in \llbracket 1, r \rrbracket$) as an $(r-1)$ -dimensional distribution. If the m^i can be chosen in such a way as to make all joint distributions equal, then the Markov kernels in the sequence $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$ are called compatible. And if no choice of $(m^i)_{i \in \llbracket 1, r \rrbracket}$ can make all joint distributions equal, we have to look for a ‘‘compromise’’ between the Markov kernels.

Remark (Producing incompatibility is easy). Take $r = 2$ and $\Omega_1 = \Omega_2 = \mathbb{R}$ and endow \mathbb{R} with the Borel sigma-algebra. Assume that for every $t \in \mathbb{R}$, $\pi_1(\cdot|t)$ and $\pi_2(\cdot|t)$ are absolutely continuous with respect to the Lebesgue measure and denote by $f_1(\cdot|t)$ and $f_2(\cdot|t)$ their respective density functions. A necessary condition [Arnold et al., 2001] for the compatibility of π_1 and π_2 is the existence of two mappings u and v defined on \mathbb{R} such that for almost every real numbers x and t (in the sense of the Lebesgue measure),

$$\frac{f_1(x|t)}{f_2(t|x)} = u(t)v(x). \quad (4.1)$$

Definition 4.2. Let ϕ be a probability distribution on \mathcal{A} . For every $i \in \llbracket 1, r \rrbracket$, denote by ϕ_{-i} the probability distribution on $\otimes_{j \neq i} \mathcal{A}_j$ defined as follows. For every set S_{-i} that can be decomposed as $S_{-i} = \times_{j \neq i} S_j$ (with $S_j \in \mathcal{A}_j$ for every $j \neq i$),

$$\phi_{-i}(S_{-i}) = \phi\left(\times_{j < i} S_j \times \Omega_i \times \times_{k > i} S_k\right). \quad (4.2)$$

ϕ_{-i} is called the i -th $(r-1)$ -marginal distribution of ϕ .

Remark. The above definition is valid because any probability distribution on \mathcal{A} can be characterized by its values on ‘‘rectangles’’ $\times_{i=1}^r S_i$ (where for every $i \in \llbracket 1, r \rrbracket$, $S_i \in \mathcal{A}_i$).

For every $i \in \llbracket 1, r \rrbracket$ and every probability distribution m^i on $\otimes_{j \neq i} \mathcal{A}_j$, denote by $\pi_i m^i$ the distribution on \mathcal{A} defined as follows. For every $i \in \llbracket 1, r \rrbracket$, for every set $S_{<i} \in \otimes_{j < i} \mathcal{A}_j$, every set $S_{>i} \in \otimes_{k > i} \mathcal{A}_k$ and every set $S_i \in \mathcal{A}_i$,

$$\pi_i m^i(S_{<i} \times S_i \times S_{>i}) = \int_{S_{<i} \times S_{>i}} \pi_i(S_i | \omega_{-i}) dm^i(\omega_{-i}). \quad (4.3)$$

Naturally, for $i = 1$ (resp. $i = r$), remove $S_{<i}$ (resp. $S_{>i}$) from the formula above. In the following, do this kind of operation when $i = 1$ or $i = r$.

Notice that for every $i \in \llbracket 1, r \rrbracket$, m^i is the i -th $(r-1)$ -marginal distribution of $\pi_i m^i$:

$$(\pi_i m^i)_{-i} = m^i. \quad (4.4)$$

If there exists a sequence of $(r-1)$ -dimensional distributions $(m^i)_{i \in \llbracket 1, r \rrbracket}$ such that all distributions $\pi_i m^i$ are equal, then the Markov kernels $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$ are compatible. If no such sequence $(m^i)_{i \in \llbracket 1, r \rrbracket}$ exists, then we wish to find a sequence $(m^i)_{i \in \llbracket 1, r \rrbracket}$ that makes the $\pi_i m^i$ share some ‘‘common ground’’. The following definition expresses this constraint formally.

Definition 4.3. A sequence of $(r - 1)$ -dimensional distributions $(m^i)_{i \in \llbracket 1, r \rrbracket}$ (each m^i being a probability distribution on $\bigotimes_{j \neq i} \mathcal{A}_j$) is said to be compatible with the sequence of Markov kernels $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$ if for every $i \in \llbracket 1, r \rrbracket$

$$m^i = \frac{1}{r} \sum_{j=1}^r (\pi_j m^j)_{-i}. \quad (4.5)$$

So the “common ground” we require for the sequence of distributions $(\pi_i m^i)_{i \in \llbracket 1, r \rrbracket}$ is that their $(r - 1)$ -marginal distributions should be the same on average. Other constraints would have been possible, and we discuss some of them in Section 4.3 below.

The definition of a compromise follows from this new definition of compatibility.

Definition 4.4. A probability distribution P on \mathcal{A} is called a compromise between the Markov kernels $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$ if these two conditions are verified:

1. for every $i \in \llbracket 1, r \rrbracket$, $\pi_i P_{-i}$ is absolutely continuous with respect to P ;
2. the sequence $(P_{-i})_{i \in \llbracket 1, r \rrbracket}$ of P 's $(r - 1)$ -marginal distributions is compatible with the sequence of Markov kernels $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$.

In the definition of a compromise, the first condition exists to give meaning to the definition of an optimal compromise below. It is reasonable on its own though: a compromise should not deem events impossible if they are considered possible by the Markov kernels.

Definition 4.5. Let λ be a positive measure on \mathcal{A} . Let P be a compromise between the sequence of Markov kernels $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$ that is absolutely continuous with respect to λ . P is called an optimal compromise with respect to λ between the sequence of Markov kernels $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$ if it minimizes the functional E_λ over all compromises between $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$ that are absolutely continuous with respect to λ . E_λ is defined by:

$$E_\lambda(P) = \sum_{i=1}^r \int_{\mathcal{A}} \left[\frac{d(\pi_i P_{-i})}{d\lambda}(\omega) - \frac{dP}{d\lambda}(\omega) \right]^2 d\lambda(\omega). \quad (4.6)$$

Remark. The set of all compromises between $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$ is convex, as is the subset of all compromises absolutely continuous with respect to λ .

If the Markov kernels $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$ are compatible and there exists a joint distribution π on \mathcal{A} that agrees with them all, then for every positive measure λ on \mathcal{A} such that π is absolutely continuous with respect to λ , $E_\lambda(\pi) = 0$ and π is an optimal compromise with respect to λ .

Even though Definition 4.5 makes it seem like the notion of optimal compromise is tied to a reference measure λ , it turns out that in many situations there exists a compromise that is optimal with respect to all possible reference measures.

Deriving the optimal compromise

The notion of Gibbs compromise is central to this theory of compromises between incompatible Markov kernels.

Definition 4.6. A probability distribution P_G on \mathcal{A} is called a Gibbs compromise between the sequence of Markov kernels $(\pi_i)_{i \in [1, r]}$ if it satisfies:

$$P_G = \frac{1}{r} \sum_{i=1}^r \pi_i(P_G)_{-i}. \quad (4.7)$$

Proposition 4.7. A Gibbs compromise between the sequence of Markov kernels $(\pi_i)_{i \in [1, r]}$ is also a compromise between this sequence of Markov kernels in the sense of Definition 4.4.

Proof. Let P_G be a Gibbs compromise between the sequence of Markov kernels $(\pi_i)_{i \in [1, r]}$. Then, for any measurable set A such that $P_G(A) = 0$,

$$P_G = \frac{1}{r} \sum_{i=1}^r \pi_i(P_G)_{-i}(A) = 0. \quad (4.8)$$

So for every integer $i \in [1, r]$, $\pi_i(P_G)_{-i}(A) = 0$. This fulfills the first condition of Definition 4.4.

Its second condition is also fulfilled because for every integer $i \in [1, r]$,

$$(P_G)_{-i} = \left(\frac{1}{r} \sum_{j=1}^r \pi_j(P_G)_{-j} \right)_{-i} = \frac{1}{r} \sum_{j=1}^r (\pi_j(P_G)_{-j})_{-i}. \quad (4.9)$$

□

The denomination ‘‘Gibbs compromise’’ is justified because it is a stationary distribution for the Gibbs sampler with random equiprobable scanning order.

The proposition below shows that all compromises are tied to Gibbs compromises.

Proposition 4.8. If a sequence of $(r - 1)$ -dimensional probability distributions $(m^i)_{i \in [1, r]}$ (each m^i being a probability distribution on $\bigotimes_{j \neq i} \mathcal{A}_j$) is compatible with the sequence of Markov kernels $(\pi_i)_{i \in [1, r]}$, then it is the sequence of $(r - 1)$ -dimensional distributions of a Gibbs compromise between the Markov kernels $(\pi_i)_{i \in [1, r]}$.

Proof. Define the probability distribution P on \mathcal{A} as follows:

$$P = \frac{1}{r} \sum_{i=1}^r \pi_i m^i. \quad (4.10)$$

Then for every $i \in [1, r]$ the i -th $(r - 1)$ -marginal distribution P_{-i} of P is given by

$$P_{-i} = \frac{1}{r} \sum_{j=1}^r (\pi_j m^j)_{-i} = m^i, \quad (4.11)$$

where the last equality is due to $(m^i)_{i \in [1, r]}$ being compatible with $(\pi_i)_{i \in [1, r]}$. Plugging this into Equation (4.10), we obtain that P is a Gibbs compromise. Equation (4.11) then yields the result. □

Remark (Equivalence relationship and convexity). Let us say that two compromises are equivalent if they share the same sequence of $(r - 1)$ -marginal distributions. This is obviously an equivalence relationship and every class of equivalence can be represented by a single Gibbs

compromise. Moreover, each class of equivalence is a convex subset of the set of all compromises. And for any positive measure λ on \mathcal{A} , its intersection with the set of all compromises absolutely continuous with respect to λ is also convex. Finally, the functional E_λ is convex over this intersection.

The main result follows from Proposition 4.8:

Theorem 4.9. *If there exists a unique Gibbs compromise P_G between the sequence of Markov kernels $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$, then the following statements holds:*

1. P_G is absolutely continuous with respect to any compromise between the sequence of Markov kernels $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$;
2. for any positive measure λ on \mathcal{A} such that P_G is absolutely continuous with respect to λ , P_G is the unique optimal compromise with respect to λ .

Because of these two properties, we call P_G the optimal compromise.

Proof. If the distribution P_G is the unique Gibbs compromise between the sequence of Markov kernels $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$, then Proposition 4.8 asserts that $((P_G)_{-i})_{i \in \llbracket 1, r \rrbracket}$ is the only compatible sequence of $(r-1)$ -dimensional distributions. So any compromise has the same sequence of $(r-1)$ -marginal distributions. Let P be such a compromise.

For every $i \in \llbracket 1, r \rrbracket$, $\pi_i(P_G)_{-i} = \pi_i P_{-i}$ is absolutely continuous with respect to P . So the average P_G is also absolutely continuous with respect to P .

Let λ be a positive measure on \mathcal{A} such that P is absolutely continuous with respect to λ . Because P_G and P share the same sequence of $(r-1)$ -marginal distributions, for λ -almost any $\omega \in \Omega$,

$$\frac{d(\pi_i P_{-i})}{d\lambda}(\omega) = \frac{d(\pi_i (P_G)_{-i})}{d\lambda}(\omega). \quad (4.12)$$

Moreover, Equation (4.10) implies that for λ -almost any $\omega \in \Omega$, $\frac{dP_G}{d\lambda}(\omega)$ is the arithmetic average between the $\frac{d(\pi_i (P_G)_{-i})}{d\lambda}(\omega)$ ($i \in \llbracket 1, r \rrbracket$), so it minimizes the mean squared error. Together with Equation (4.12), this implies that for λ -almost any $\omega \in \Omega$,

$$\sum_{i=1}^r \left[\frac{d(\pi_i (P_G)_{-i})}{d\lambda}(\omega) - \frac{dP_G}{d\lambda}(\omega) \right]^2 \leq \sum_{i=1}^r \left[\frac{d(\pi_i P_{-i})}{d\lambda}(\omega) - \frac{dP}{d\lambda}(\omega) \right]^2. \quad (4.13)$$

Consequently, $E_\lambda(P_G) \leq E_\lambda(P)$. Moreover, if $P \neq P_G$, then there exists $S \in \mathcal{A}$ such that $\lambda(S) > 0$ and for every $\omega \in S$,

$$\forall i \in \llbracket 1, r \rrbracket \quad \frac{d(\pi_i P_{-i})}{d\lambda}(\omega) = \frac{d(\pi_i (P_G)_{-i})}{d\lambda}(\omega) \quad \text{and} \quad \frac{dP}{d\lambda}(\omega) \neq \frac{d(P_G)}{d\lambda}(\omega). \quad (4.14)$$

So for every $\omega \in S$, Equation (4.13) is a strict inequality and thus $E_\lambda(P_G) < E_\lambda(P)$. P_G is therefore the unique optimal compromise with respect to λ . \square

Remark (Equivalence class and convexity – continued). A similar proof can be used to show the following results. In cases where there are several different Gibbs compromises, each Gibbs compromise is absolutely continuous with respect to any equivalent compromise. Now let \mathcal{C} be an equivalence class and let $\pi_{\mathcal{C}}$ be the unique Gibbs compromise in this class of equivalence. Let λ be a positive measure on \mathcal{A} such that $\pi_{\mathcal{C}}$ is absolutely continuous with respect to λ . $\pi_{\mathcal{C}}$ (uniquely) minimizes E_λ over the intersection of \mathcal{C} and the set of all compromises absolutely

continuous with respect to λ . This implies that for any positive measure λ on \mathcal{A} , any optimal compromise with respect to λ is a Gibbs compromise. Finally, note that the set of all Gibbs compromises is convex (however there is no reason E_λ should be convex over this set!).

Theorem 4.9 has important practical implications. It opens the possibility of using Gibbs sampling to find the optimal compromise between incompatible Markov kernels and justifies using PIGS.

The next result shows that, under the conditions of Theorem 4.9, the optimal compromise remains the same for a fairly large class of reparametrizations. This result is key to the application of PIGS in an Objective Bayesian framework, where some degree of invariance by reparametrization of priors and posteriors is usually expected.

For every $i \in \llbracket 1, r \rrbracket$, let $(\tilde{\Omega}_i, \tilde{\mathcal{A}}_i)$ be a measurable space and let f_i be bijective measurable mapping $\Omega_i \rightarrow \tilde{\Omega}_i$ whose inverse f_i^{-1} is also measurable. Define $f = (f_1, \dots, f_r) : \times_{i \in \llbracket 1, r \rrbracket} \Omega_i \rightarrow \times_{i \in \llbracket 1, r \rrbracket} \tilde{\Omega}_i$ and for every $i \in \llbracket 1, r \rrbracket$ $f_{-i} = (f_1, \dots, f_{i-1}, f_{i+1}, \dots, f_r) : \times_{j \neq i} \Omega_j \rightarrow \times_{j \neq i} \tilde{\Omega}_j$.

Also let $\tilde{\pi}_i$ be the Markov kernel $(\times_{j \neq i} \tilde{\Omega}_j) \times \tilde{\mathcal{A}}_i \rightarrow [0, 1]$ such that for every $\omega_{-i} \in \times_{j \neq i} \Omega_j$ and every $S_i \in \mathcal{A}_i$, $\tilde{\pi}_i(f_i(S_i) | f_{-i}(\omega_{-i})) = \pi_i(S_i | \omega_{-i})$.

Proposition 4.10. *Assume there exists a unique Gibbs compromise P_G between the sequence of Markov kernels $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$. Then the push-forward measure of P_G by f $\tilde{P}_G := P_G * f$ is the unique Gibbs compromise between the sequence of Markov kernels $(\tilde{\pi}_i)_{i \in \llbracket 1, r \rrbracket}$.*

Proof. For every $i \in \llbracket 1, r \rrbracket$, for every $\tilde{S}_i \in \tilde{\mathcal{A}}_i$,

$$\begin{aligned}
\tilde{P}_G \left(\times_{i \in \llbracket 1, r \rrbracket} \tilde{S}_i \right) &= P_G \left(f^{-1} \left(\times_{i \in \llbracket 1, r \rrbracket} \tilde{S}_i \right) \right) \\
&= P_G \left(\times_{i \in \llbracket 1, r \rrbracket} f_i^{-1}(\tilde{S}_i) \right) \\
&= \frac{1}{r} \sum_{i=1}^r \int_{\times_{j \neq i} f_j^{-1}(\tilde{S}_j)} \pi_i(f_i^{-1}(\tilde{S}_i) | \omega_{-i}) d\{(P_G)_{-i}\}(\omega_{-i}) \\
&= \frac{1}{r} \sum_{i=1}^r \int_{\times_{j \neq i} f_j^{-1}(\tilde{S}_j)} \tilde{\pi}_i(\tilde{S}_i | f_{-i}(\omega_{-i})) d\{(P_G)_{-i}\}(\omega_{-i}) \\
&= \frac{1}{r} \sum_{i=1}^r \int_{\times_{j \neq i} \tilde{S}_j} \tilde{\pi}_i(\tilde{S}_i | \tilde{\omega}_{-i}) d\{(P_G)_{-i} * f_{-i}\}(\tilde{\omega}_{-i}). \tag{4.15}
\end{aligned}$$

Now, for every $i \in \llbracket 1, r \rrbracket$, for every $\tilde{T}_i \in \tilde{\mathcal{A}}_i$,

$$\begin{aligned}
(P_G)_{-i} * f_{-i} \left(\times_{j \neq i} \tilde{T}_j \right) &= (P_G)_{-i} \left(\times_{j \neq i} f_j^{-1}(\tilde{T}_j) \right) \\
&= P_G \left(\times_{j < i} f_j^{-1}(\tilde{T}_j) \times f_i^{-1}(\tilde{\Omega}_i) \times \times_{k > i} f_k^{-1}(\tilde{T}_k) \right) \\
&= P_G * f \left(\times_{j < i} \tilde{T}_j \times \tilde{\Omega}_i \times \times_{k > i} \tilde{T}_k \right) \\
&= (P_G * f)_{-i} \left(\times_{j \neq i} \tilde{T}_j \right). \tag{4.16}
\end{aligned}$$

So $(P_G)_{-i} * f_{-i} = (P_G * f)_{-i} = (\tilde{P}_G)_{-i}$. Then, returning to Equation 4.15,

$$\tilde{P}_G \left(\times_{i \in \llbracket 1, r \rrbracket} \tilde{S}_i \right) = \frac{1}{r} \sum_{i=1}^r \int_{\times_{j \neq i} \tilde{S}_j} \tilde{\pi}_i(\tilde{S}_i | \tilde{\omega}_{-i}) d \left\{ (\tilde{P}_G)_{-i} \right\} (\tilde{\omega}_{-i}). \quad (4.17)$$

Finally, we obtain

$$\tilde{P}_G = \frac{1}{r} \sum_{i=1}^r \tilde{\pi}_i (\tilde{P}_G)_{-i}. \quad (4.18)$$

\tilde{P}_G is therefore a Gibbs compromise between the sequence of Markov kernels $(\tilde{\pi}_i)_{i \in \llbracket 1, r \rrbracket}$. We now prove its uniqueness. For every $i \in \llbracket 1, r \rrbracket$, f_i is bijective and f_i^{-1} is measurable. So for any Gibbs compromise \tilde{Q}_G between the sequence of Markov kernels $(\tilde{\pi}_i)_{i \in \llbracket 1, r \rrbracket}$, $\tilde{Q}_G * f^{-1}$ is a Gibbs compromise between the sequence of Markov kernels $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$. Given P_G is the unique Gibbs compromise between the Markov kernels in this sequence, $\tilde{Q}_G * f^{-1} = P_G$, so $\tilde{Q}_G = \tilde{Q}_G * f^{-1} * f = P_G * f = \tilde{P}_G$. □

4.3 Discussion of the notion of compromise

This section is devoted to discussing the merits, in our opinion, of Definitions 4.3 and 4.4. The first subsection shows that strengthening their requirements often makes the set of all compromises empty. Conversely, the second subsection provides examples showing that relaxing their requirements leads to undesirable behavior in compromises.

Stronger definitions of compromises are not possible

While the definition of the *optimal compromise* is straightforward, as it involves minimizing some measure of distance between the “targeted” conditionals and the conditionals of the compromise, the definition of a *compromise* may seem arbitrary. To motivate this definition, let us focus on the two-dimensional case.

Suppose that $r = 2$ and that π_1 and π_2 are incompatible. This means there exists no joint distribution π which agrees with both Markov kernels. This being the case, it seems sensible to weaken the definition of compatibility by applying it to the “marginals” instead of the “joint” distribution. The following definition makes this idea precise.

Definition 4.11. *A pair of probability distributions m^1 (resp. m^2) on \mathcal{A}_2 (resp. \mathcal{A}_1) is compatible with the pair of Markov kernels π_1 and π_2 if the distributions $\pi_1 m^1$ and $\pi_2 m^2$ verify*

$$(\pi_1 m^1)_{-2} = m^2 \text{ and } (\pi_2 m^2)_{-1} = m^1. \quad (4.19)$$

While this definition may seem more restrictive at first glance than Definition 4.3, both definitions are in fact equivalent when applied to a pair of Markov kernels, because $(r - 1)$ -dimensional distributions are simply 1-dimensional distributions in this case. Indeed, following directly from Definition 4.3, we have this result which holds for any r :

Proposition 4.12. *If a sequence of $(r - 1)$ -dimensional probability distributions $(m^i)_{i \in \llbracket 1, r \rrbracket}$ (each m^i being a probability distribution on $\otimes_{j \neq i} \mathcal{A}_j$) is compatible (in the sense of Definition*

4.3) with the sequence of Markov kernels $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$, then all joint distributions in the sequence $(\pi_i m^i)_{i \in \llbracket 1, r \rrbracket}$ share the same marginals, that is $\forall i, j, k \in \llbracket 1, r \rrbracket$,

$$\forall S_k \in \mathcal{A}_k, \quad \pi_i m^i \left(\prod_{k' < k} \Omega_{k'} \times S_k \times \prod_{k'' > k} \Omega_{k''} \right) = \pi_j m^j \left(\prod_{k' < k} \Omega_{k'} \times S_k \times \prod_{k'' > k} \Omega_{k''} \right). \quad (4.20)$$

Now let us consider the three-dimensional case ($r = 3$). Because the aim of this section is merely to motivate the definitions of compromises and optimal compromises, there is no need for the discussion to be fully general. Let us therefore restrict the discussion to an important particular case. Assume that Ω_1 , Ω_2 and Ω_3 are finite sets and that \mathcal{A}_1 , \mathcal{A}_2 and \mathcal{A}_3 are respectively the sets of all their subsets. This has an important consequence: any mapping from a subset of Ω to a subset of Ω is measurable.

Also assume that the Markov kernels π_1 , π_2 and π_3 are positive mappings. For π_1 , this means that for every $(\omega_1, \omega_2, \omega_3) \in \Omega_1 \times \Omega_2 \times \Omega_3$, $\pi_1(\{\omega_1\} | \omega_2, \omega_3) > 0$.

If we consider ω_3 known, then the situation is reduced to the two-dimensional case. Because π_1 and π_2 are positive mappings and both Ω_1 and Ω_2 are finite sets, Markov chain theory ensures there exists a unique Gibbs compromise $P(\cdot | \omega_3)$. For every $(\omega_1, \omega_2) \in \Omega_1 \times \Omega_2$,

$$P(\{\omega_1\} \times \{\omega_2\} | \omega_3) = \frac{1}{2} \pi_1(\{\omega_1\} | \omega_2, \omega_3) P(\Omega_1 \times \{\omega_2\} | \omega_3) + \frac{1}{2} \pi_1(\omega_2 | \omega_1, \omega_3) P(\{\omega_1\} \times \Omega_2 | \omega_3). \quad (4.21)$$

Thanks to Theorem 4.9, $P(\cdot | \omega_3)$ is the optimal compromise. Moreover, notice that Equation (4.21) defines a Markov kernel on $\Omega_3 \times (\mathcal{A}_1 \otimes \mathcal{A}_2)$.

We may similarly derive Markov kernels $Q : \Omega_1 \times (\mathcal{A}_2 \otimes \mathcal{A}_3)$ and $R : \Omega_2 \times (\mathcal{A}_1 \otimes \mathcal{A}_3)$.

Once again, because π_1 , π_2 and π_3 are positive mappings, it follows from Markov chain theory that P , Q and R are also positive mappings. We now show it using only elementary arguments, because these arguments will be useful again later.

Assume P is not a positive mapping. Then there exists $(\omega_1^{(0)}, \omega_2^{(0)}, \omega_3^{(0)}) \in \Omega_1 \times \Omega_2 \times \Omega_3$ such that $P(\{\omega_1^{(0)}\} \times \{\omega_2^{(0)}\} | \omega_3^{(0)}) = 0$. Equation (4.21) then implies that $P(\Omega_1 \times \{\omega_2^{(0)}\} | \omega_3^{(0)}) = 0$. So for every $\omega_1 \in \Omega_1$, $P(\{\omega_1\} \times \{\omega_2^{(0)}\} | \omega_3^{(0)}) = 0$. But then Equation (4.21) implies that $P(\{\omega_1\} \times \Omega_2 | \omega_3^{(0)}) = 0$. Since this holds for every $\omega_1 \in \Omega_1$, $P(\Omega_1 \times \Omega_2 | \omega_3^{(0)}) = 0$, which is absurd since $P(\cdot | \omega_3^{(0)})$ is supposed to be a probability distribution. So P is a positive mapping.

Ideally, we would wish to define the optimal compromise between π_1 , π_2 and π_3 as the joint distribution ϕ such that for every $(\omega_1, \omega_2, \omega_3) \in \Omega_1 \times \Omega_2 \times \Omega_3$,

$$\phi(\{\omega_1\} \times \{\omega_2\} \times \{\omega_3\}) = P(\{\omega_1\} \times \{\omega_2\} | \omega_3) \phi(\Omega_1 \times \Omega_2 \times \{\omega_3\}) \quad (4.22)$$

$$= Q(\{\omega_2\} \times \{\omega_3\} | \omega_1) \phi(\{\omega_1\} \times \Omega_2 \times \Omega_3) \quad (4.23)$$

$$= R(\{\omega_1\} \times \{\omega_3\} | \omega_2) \phi(\Omega_2 \times \{\omega_2\} \times \Omega_3). \quad (4.24)$$

Unfortunately, the existence of such an ‘‘optimal compromise’’ ϕ implies that π_1 , π_2 and π_3 are compatible. First, one can show that for every $(\omega_1, \omega_2, \omega_3) \in \Omega_1 \times \Omega_2 \times \Omega_3$, $\phi(\{\omega_1\} \times$

$\{\omega_2\} \times \{\omega_3\} > 0$. The arguments are similar to those used above to show that P is a positive mapping. Rewriting Equations (4.22) and (4.24) yields

$$\begin{aligned} & \frac{\phi(\{\omega_1\} \times \{\omega_2\} \times \{\omega_3\})}{\phi(\Omega_1 \times \{\omega_2\} \times \{\omega_3\})} \\ &= \frac{P(\{\omega_1\} \times \{\omega_2\} | \omega_3)}{P(\Omega_1 \times \{\omega_2\} | \omega_3)} = \frac{1}{2} \pi_1(\{\omega_1\} | \omega_2, \omega_3) + \frac{1}{2} \pi_2(\{\omega_2\} | \omega_1, \omega_3) \frac{P(\{\omega_1\} \times \Omega_2 | \omega_3)}{P(\Omega_1 \times \{\omega_2\} | \omega_3)} \end{aligned} \quad (4.25)$$

$$= \frac{R(\{\omega_1\} \times \{\omega_3\} | \omega_2)}{R(\Omega_1 \times \{\omega_2\} | \omega_3)} = \frac{1}{2} \pi_1(\{\omega_1\} | \omega_2, \omega_3) + \frac{1}{2} \pi_3(\{\omega_3\} | \omega_1, \omega_2) \frac{R(\{\omega_1\} \times \Omega_3 | \omega_2)}{R(\Omega_1 \times \{\omega_3\} | \omega_2)}. \quad (4.26)$$

Now combine Equations (4.25) and (4.26):

$$\frac{1}{2} \pi_2(\{\omega_2\} | \omega_1, \omega_3) \frac{\phi(\{\omega_1\} \times \Omega_2 \times \{\omega_3\})}{\phi(\Omega_1 \times \{\omega_2\} \times \{\omega_3\})} = \frac{1}{2} \pi_3(\{\omega_3\} | \omega_1, \omega_2) \frac{\phi(\{\omega_1\} \times \{\omega_2\} \times \Omega_3)}{\phi(\Omega_1 \times \{\omega_2\} \times \{\omega_3\})}. \quad (4.27)$$

As this holds for every $(\omega_1, \omega_2, \omega_3) \in \Omega_1 \times \Omega_2 \times \Omega_3$, it implies that $\pi_2 \phi_{-2} = \pi_3 \phi_{-3}$. A similar proof then shows that $\pi_3 \phi_{-3} = \pi_1 \phi_{-1}$. This means that if an “optimal compromise” ϕ exists, then π_1 , π_2 and π_3 are compatible and no compromise was needed.

Similarly to what was done in the two-dimensional case, we avoid this difficulty by weakening the compatibility requirements: we no longer require $P(\cdot | \omega_3)$ to be the optimal compromise between π_1 and π_2 for every $\omega_3 \in \Omega_3$ (as expressed by Equation (4.21)), but only on average over $\omega_3 \in \Omega_3$. So a compromise ϕ should still verify Equation (4.22), but it would only need to verify this weakened version of Equation (4.21) for every $(\omega_1, \omega_2) \in \Omega_1 \times \Omega_2$:

$$\begin{aligned} & \sum_{\omega_3 \in \Omega_3} P(\{\omega_1\} \times \{\omega_2\} | \omega_3) \phi(\Omega_1 \times \Omega_2 \times \{\omega_3\}) \\ &= \sum_{\omega_3 \in \Omega_3} \left[\frac{1}{2} \pi_1(\{\omega_1\} | \omega_2, \omega_3) P(\Omega_1 \times \{\omega_2\} | \omega_3) + \frac{1}{2} \pi_2(\{\omega_2\} | \omega_1, \omega_3) P(\{\omega_1\} \times \Omega_2 | \omega_3) \right] \times \\ & \quad \phi(\Omega_1 \times \Omega_2 \times \{\omega_3\}). \end{aligned} \quad (4.28)$$

Because Equation (4.22) is still expected to hold, Equation (4.28) is equivalent to

$$\begin{aligned} & \phi_{-3}(\{\omega_1\} \times \{\omega_2\}) \\ &= \frac{1}{2} \sum_{\omega_3 \in \Omega_3} \pi_1(\{\omega_1\} | \omega_2, \omega_3) \phi_{-1}(\{\omega_2\} \times \{\omega_3\}) + \frac{1}{2} \sum_{\omega_3 \in \Omega_3} \pi_2(\{\omega_2\} | \omega_1, \omega_3) \phi_{-2}(\{\omega_1\} \times \{\omega_3\}). \end{aligned} \quad (4.29)$$

So the requirement boils down to $2\phi_{-3} = (\pi_1 \phi_{-1})_{-3} + (\pi_2 \phi_{-2})_{-3}$. Of course, we symmetrically require $2\phi_{-1} = (\pi_2 \phi_{-2})_{-1} + (\pi_3 \phi_{-3})_{-1}$ and $2\phi_{-2} = (\pi_1 \phi_{-1})_{-2} + (\pi_3 \phi_{-3})_{-2}$ as well.

To sum this part of the discussion up, the necessity of weakening our “ideal” requirements for “optimal compatibility” made us downgrade from a requirement about a “joint” 3-dimensional distribution ϕ to requirements about its $(3 - 1)$ -marginal distributions ϕ_{-1} , ϕ_{-2} and ϕ_{-3} . Definition 4.3 is just another formulation of these requirements, which are taken to define a compatible sequence of $(r - 1)$ -dimensional distributions.

Proposition 4.8 shows that in cases where there exists a unique Gibbs compromise, even with this weakened set of requirements for compatibility, there exists only one compatible sequence

of $(r - 1)$ -dimensional distributions, so we may not strengthen it if there is to be any solution. Indeed, in cases where no Gibbs compromise exists, no compatible set of $(r - 1)$ -dimensional distributions exists either!

Weaker definitions of compromises are inconvenient

As was shown in the previous subsection, the requirements given by Definition 4.3 for the compatibility of a sequence of $(r - 1)$ -dimensional distributions with a given sequence of Markov kernels cannot be strengthened. The following shows they cannot be weakened either.

Why we need some notion of compatibility: example in 2 dimensions.

Why bother with the compatibility of $(r - 1)$ -dimensional distributions and not simply minimize the functional E_λ of Definition 4.5 over all distributions absolutely continuous with respect to λ ? The following 2-dimensional example shows that doing so yields unsatisfactory results.

Consider the following situation: $\Omega_1 = \Omega_2 = \{0, 1\}$ and $\mathcal{A}_1 = \mathcal{A}_2 = \{\emptyset, \{0\}, \{1\}, \{0, 1\}\}$. Let X_1 (resp. X_2) be the identity function on Ω_1 (resp. Ω_2). Both X_1 and X_2 are measurable functions (and thus random variables when $\mathcal{A}_1 \otimes \mathcal{A}_2$ is endowed with a probability measure). Define the following Markov kernels:

$$\pi_1(X_1 = 1|\omega_2) = \mathbf{1}_{\{\omega_2=0\}} + 1/2 \mathbf{1}_{\{\omega_2=1\}}; \quad (4.30)$$

$$\pi_2(X_2 = 1|\omega_1) = 1/2 \mathbf{1}_{\{\omega_1=0\}} + \mathbf{1}_{\{\omega_1=1\}}. \quad (4.31)$$

In this situation, let us define the reference measure λ as the counting measure.

Denote by π_C the optimal compromise (maximizing E_λ over all compromises) and π_E the distribution on $\mathcal{A}_1 \otimes \mathcal{A}_2$ that minimizes E_λ over all distributions on $\mathcal{A}_1 \otimes \mathcal{A}_2$. We have $E(\pi_C) = 2/25 > E(\pi_E) = 1/15$ and

$$\begin{aligned} \pi_C(\{\omega_1\} \times \{\omega_2\}) &= 1/10 (\mathbf{1}_{\{(\omega_1, \omega_2)=(0,0)\}} + \mathbf{1}_{\{(\omega_1, \omega_2)=(1,0)\}} \\ &\quad + 3 \mathbf{1}_{\{(\omega_1, \omega_2)=(0,1)\}} + 5 \mathbf{1}_{\{(\omega_1, \omega_2)=(1,1)\}}); \end{aligned} \quad (4.32)$$

$$\begin{aligned} \pi_E(\{\omega_1\} \times \{\omega_2\}) &= 1/30 (3 \mathbf{1}_{\{(\omega_1, \omega_2)=(0,0)\}} + \mathbf{1}_{\{(\omega_1, \omega_2)=(1,0)\}} \\ &\quad + 11 \mathbf{1}_{\{(\omega_1, \omega_2)=(0,1)\}} + 15 \mathbf{1}_{\{(\omega_1, \omega_2)=(1,1)\}}). \end{aligned} \quad (4.33)$$

As π_C is the unique Gibbs compromise between π_1 and π_2 , its marginals are the only compatible marginals: $(\pi_C)_{-1}(X_2 = 1) = 4/5$ and $(\pi_C)_{-2}(X_1 = 1) = 3/5$. The marginals of π_E are noticeably different: $(\pi_E)_{-1}(X_2 = 1) = 13/15$ and $(\pi_E)_{-2}(X_1 = 1) = 8/15$.

Observe that according to π_1 , $X_2 = 0$ implies $X_1 = 1$ but that according to π_2 , $X_1 = 1$ implies that $X_2 = 1 \neq 0$. This discrepancy is a major source of incompatibility between the two Markov kernels. So, as π_E makes both $X_1 = 1$ and $X_2 = 0$ less likely than π_C , it “ignores the inconsistent parts” of π_1 and π_2 to some extent. Therefore, if the marginals are not set in advance (say, by imposing compatibility with the conditionals in the sense of Definition 4.3), one may “cheat” by having the marginals disadvantage inconvenient values for the parameters.

Why the compatibility requirements cannot be weakened: example in 3 dimensions.

In the two-dimensional case, because of Proposition 4.12, Definitions 4.11 and 4.3 give the same meaning to the concept of compatibility of marginals, so Definition 4.3 may be thought of as a generalization of Definition 4.11 to cases with more than two dimensions. However, another generalization of 4.11 is possible. To avoid confusion, this other generalization will be called *weak compatibility*. In the following, the sequence of Markov kernels $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$ is defined as in Section 4.2.

Definition 4.13. *A sequence of $(r - 1)$ -dimensional distributions $(m^i)_{i \in \llbracket 1, r \rrbracket}$ is weakly compatible with a sequence of Markov kernels $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$ if Equation (4.20) holds.*

Proposition 4.12 means compatibility in the sense of Definition 4.3 implies weak compatibility in the sense of Definition 4.13, hence its denomination as “weak”.

Using the concept of weak compatibility of a sequence of $(r - 1)$ -marginal distributions, we define weak compromises and the optimal weak compromise as analogues to compromises and optimal compromises respectively.

Definition 4.14. *A probability distribution P on $\bigotimes_{i \in \llbracket 1, r \rrbracket} \mathcal{A}_i$ is called a weak compromise between the Markov kernels $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$ if these two conditions are verified:*

1. *for every $i \in \llbracket 1, r \rrbracket$, $\pi_i P_{-i}$ is absolutely continuous with respect to P ;*
2. *the sequence $(P_{-i})_{i \in \llbracket 1, r \rrbracket}$ of P 's $(r - 1)$ -marginal distributions is weakly compatible with $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$.*

Definition 4.15. *Let λ be a positive measure on \mathcal{A} . Let P be a weak compromise between the sequence of Markov kernels $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$ that is absolutely continuous with respect to λ . P is called an optimal weak compromise with respect to λ between the sequence of Markov kernels $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$ if it minimizes the functional E_λ over all compromises between $(\pi_i)_{i \in \llbracket 1, r \rrbracket}$ that are absolutely continuous with respect to λ . E_λ is defined by Equation (4.6).*

Because, for any positive measure λ on \mathcal{A} , the set of all weak compromises absolutely continuous with respect to λ includes the set of all compromises absolutely continuous with respect to λ , an optimal weak compromise with respect to λ makes the functional E_λ no greater than an optimal compromise with respect to λ . However, as shown in the following example with $r = 3$, optimal weak compromises may have undesirable behavior.

Assume $\Omega_1 = \Omega_2 = \Omega_3 = \{0, 1\}$ and $\mathcal{A}_1 = \mathcal{A}_2 = \mathcal{A}_3 = \{\emptyset, \{0\}, \{1\}, \{0, 1\}\}$. Let X_1 (resp. X_2 , X_3) be the identity function on Ω_1 (resp. Ω_2 , Ω_3).

Consider the following Markov kernels. For every $(\omega_1, \omega_2, \omega_3) \in \Omega_1 \times \Omega_2 \times \Omega_3$

$$\pi_1(X_1 = 1 | \omega_2, \omega_3) = 1/2; \tag{4.34}$$

$$\pi_2(X_2 = 1 | \omega_1, \omega_3) = 1/2; \tag{4.35}$$

$$\pi_3(X_3 = 1 | \omega_1, \omega_2) = \mathbf{1}_{\{\omega_1 = \omega_2\}} + 1/2 \mathbf{1}_{\{\omega_1 \neq \omega_2\}}. \tag{4.36}$$

Notice that, provided ω_3 is known, π_1 and π_2 are compatible Markov kernels. The unique probability distribution on $\mathcal{A}_1 \otimes \mathcal{A}_2$ that fits both π_1 and π_2 (conditional to ω_3) verifies for every $(\omega_1, \omega_2) \in \Omega_1 \times \Omega_2$

$$P(\{\omega_1\} \times \{\omega_2\} | \omega_3) = 1/4. \quad (4.37)$$

Thus, any joint distribution fitting the Markov kernel P would make X_1 , X_2 and X_3 mutually independent. Unfortunately, no such joint distribution could fit π_3 , but we may expect compromises between π_1 , π_2 and π_3 to retain the independence of X_1 and X_2 .

Let λ be the counting measure on \mathcal{A} .

Denote by π_C the (unique) optimal compromise between π_1 , π_2 and π_3 with respect to λ : we have $E_\lambda(\pi_C) = 1/48 \approx 0.021$. For every $(\omega_1, \omega_2, \omega_3) \in \Omega_1 \times \Omega_2 \times \Omega_3$

$$\begin{aligned} \pi_C(\{\omega_1\} \times \{\omega_2\} \times \{\omega_3\}) &= \frac{1}{24} \mathbf{1}_{\{\omega_1=\omega_2, \omega_3=0\}} + \frac{1}{12} \mathbf{1}_{\{\omega_1 \neq \omega_2, \omega_3=0\}} \\ &\quad + \frac{5}{24} \mathbf{1}_{\{\omega_1=\omega_2, \omega_3=1\}} + \frac{1}{6} \mathbf{1}_{\{\omega_1 \neq \omega_2, \omega_3=1\}}. \end{aligned} \quad (4.38)$$

Notably, its third 2-marginal distribution $(\pi_C)_{-3}$ verifies for every $(\omega_1, \omega_2) \in \Omega_1 \times \Omega_2$

$$(\pi_C)_{-3}(\{\omega_1\} \times \{\omega_2\}) = 1/4. \quad (4.39)$$

So, as expected, π_C retains the independence of X_1 and X_2 . Because π_C is the unique Gibbs compromise between π_1 , π_2 and π_3 , Proposition 4.8 implies any other compromise between π_1 , π_2 and π_3 also retains this property.

Let us now consider an optimal weak compromise π_W between π_1 , π_2 and π_3 with respect to λ . Numerical computation gives us the following approximation, with $E_\lambda(\pi_W) \approx 0.019$. For every $(\omega_1, \omega_2, \omega_3) \in \Omega_1 \times \Omega_2 \times \Omega_3$,

$$\begin{aligned} \pi_W(\{\omega_1\} \times \{\omega_2\} \times \{\omega_3\}) &\approx 0.04 \mathbf{1}_{\{\omega_1=\omega_2, \omega_3=0\}} + 0.10 \mathbf{1}_{\{\omega_1 \neq \omega_2, \omega_3=0\}} \\ &\quad + 0.19 \mathbf{1}_{\{\omega_1=\omega_2, \omega_3=1\}} + 0.17 \mathbf{1}_{\{\omega_1 \neq \omega_2, \omega_3=1\}}. \end{aligned} \quad (4.40)$$

Its third 2-marginal distribution $(\pi_W)_{-3}$ is approximately

$$(\pi_W)_{-3}(\{\omega_1\} \times \{\omega_2\}) \approx 0.23 \mathbf{1}_{\{\omega_1=\omega_2\}} + 0.27 \mathbf{1}_{\{\omega_1 \neq \omega_2\}}. \quad (4.41)$$

Thus the independence of X_1 and X_2 is lost. Therefore, weak compatibility is no adequate notion of compatibility. What happened is that although $(\pi_W)_{-3}$ and $(\pi_C)_{-3}$ share the same marginals, that is

$$(\pi_W)_{-3}(\{\omega_1\} \times \Omega_2) = 1/2, \quad (4.42)$$

$$(\pi_W)_{-3}(\Omega_1 \times \{\omega_2\}) = 1/2, \quad (4.43)$$

$(\pi_W)_{-3}$ slightly disadvantages the event $X_1 = X_2$, which is where the incompatibility between π_3 and the pair (π_1, π_2) is most obvious: according to π_3 , $X_1 = X_2$ implies $X_3 = 1$, so conversely, $X_3 = 0$ should imply $X_1 \neq X_2$, when in fact π_1 and π_2 state that even given $\omega_3 = 0$, $\{X_1 \neq X_2\}$ only happens with probability 1/2. On the other hand, according to π_3 , if $X_1 \neq X_2$, then X_3 can with equal probability be 0 or 1, which matches π_1 and π_2 better.

4.4 Conclusion

In this chapter, we proposed a theory defining the notions of compromise and especially optimal compromise between potentially incompatible conditional distributions, i.e. Markov kernels. This theory is mainly derived from intuitive conceptions of what a compromise should be. In places where such conceptions were inconclusive, we relied on concrete examples to precisely determine what was acceptable and what was not in a compromise and used it to complete the definition.

One strength of this theory is that it can be applied to continuous as well as discrete probability distributions, whereas previous studies focused on the discrete, or even finite, case.

In the rest of this thesis, we focus on one particular application of this theory: deriving Objective Bayesian inference for correlation lengths in Kriging models with anisotropic correlation kernels. However, the field of potential applications is much larger. It encompasses all models that rely on conditional rather than joint probability distributions.

Part III

Application

Chapter 5

Application of the Optimal Compromise to Simple Kriging models with Matérn correlation kernels

This chapter covers the article Muré [2017], with the exception of its second section.

Abstract

The notion of optimal compromise between possibly incompatible conditional distributions allows us to perform Objective Bayesian analysis of correlation parameters in Kriging models by using univariate conditional Jeffreys-rule posterior distributions instead of the widely used multivariate Jeffreys-rule posterior. This strategy makes the full-Bayesian procedure tractable. Numerical examples show it has near-optimal frequentist performance in terms of prediction interval coverage.

Résumé

La notion de compromis optimal entre lois conditionnelles potentiellement incompatibles permet une analyse bayésienne objective des paramètres de corrélation de modèles de krigeage : nous utilisons des lois *a posteriori* de Jeffreys univariées conditionnelles plutôt que la très usitée loi *a posteriori* de Jeffreys multivariée. Cette stratégie rend la procédure pleinement bayésienne abordable. Des exemples numériques montrent que ses performances fréquentistes en termes de taux de couverture des intervalles prédictifs sont quasi-optimales.

5.1 Introduction

The present chapter is the point where all previous chapters come together.

Recalling Chapter 1, Kriging is a surrogate model used to emulate a real-valued function on a spatial domain \mathcal{D} when said function can only be evaluated on a finite subset of \mathcal{D} called “design set”. The “Kriging prediction” is the mean function of the process taken conditionally to all known values of the emulated function, i.e. the values at the points in the design set. The main advantage of the framework is its natural way of representing uncertainty about the value of the function at unobserved points [Santner et al., 2003]. Prediction does not

consist of a single value but of a complete Normal distribution. “Kriging” is the name given to the framework in the geostatistical literature [Journel and Huijbregts, 1978], but is also frequently used in the context of computer experiments and machine learning under the label “Gaussian Process regression” [Rasmussen and Williams, 2006]. In this chapter, we focus on Simple Kriging, where the Gaussian Process is assumed to be stationary with known mean, as opposed to Universal Kriging, which incorporates an unknown mean function and will be tackled in Chapter 6.

The probability distribution of a stationary Gaussian Process is characterized by a variance parameter and a correlation function (also known as “correlation kernel”) which itself depends on parameters. So one should deal with uncertainty about model parameters. The problem is “notoriously difficult”, as highlighted by Kennedy and O’Hagan [2001], because the likelihood function may often be quite flat [Li and Sudjianto, 2005]. In a Bayesian framework, this uncertainty is represented by a prior distribution on the parameters.

Recalling Chapter 2, the Objective Bayesian paradigm as explained by Berger [2006] consists in eliciting for every model a “default”, reasonable prior distribution that could be used when no explicit prior information is available. In particular, the Berger-Bernardo reference prior [Bernardo, 2005] can be algorithmically computed with minimal user intervention.

For models with a single scalar parameter, the reference prior rewards parameter values that are easily discriminated by the likelihood function. Its definition is related to the Kullback-Leibler divergence between posterior and prior [Bernardo, 2005]. For usual continuous models – essentially models where the Fisher information matrix is equal to the opposite of the expectancy of the second derivative of the log-likelihood – it coincides with the Jeffreys-rule prior.

For models with multiple parameters, the reference prior algorithm requires the user to specify an ordering on the parameters and then iteratively compute the reference prior on each parameter conditionally to all subsequent parameters. The only user input is therefore this ordering, and common sense arguments often make one more sensible than others. Of course, one could also group several parameters and treat them as one single multi-dimensional parameter, but doing so tends to produce less satisfactory inference [Berger and Bernardo, 1992].

Berger et al. [2001] were the first to derive a reference prior for the parameters of a Gaussian Process regression model –cf. Chapter 3. This model contained only one correlation parameter, however. When several correlation parameters are involved, there is no reasonable way to order them. Even if one were arbitrarily picked, computation of the prior would be analytically intractable. Several authors [Paulo, 2005, Ren et al., 2012, Kazianka and Pilz, 2012, Ren et al., 2013, Gu et al., 2018] have therefore resolved to treat all correlation parameters as a single multidimensional parameter. It is in order to avoid having to do this that we make use of Potentially Incompatible Gibbs Samplers as defined in Chapter 4.

The idea is simple: for every correlation parameter, it is possible to analytically derive the reference prior for this parameter conditionally to all others. Each of the corresponding posterior distributions can be seen as a conditional probability distribution on one correlation parameter when all others are known. These conditional distributions then serve as input to a PIGS (Potentially Incompatible Gibbs Sampler, cf. previous chapter).

Theorem 5.2 is the main result with respect to the application. First, under reasonable assumptions, the PIGS admits one single stationary probability distribution, which is the optimal compromise in the sense of the theory of Chapter 4. Second, the Markov kernel defined by the PIGS is uniformly ergodic. Since this Markov kernel is defined over an uncountable state space, the latter fact is significant. The stationary distribution, which we call the Gibbs reference posterior distribution, can be used to improve prediction of the value taken by the Gaussian process at unobserved points. Sections 5.3 and 5.4 illustrate the inferential and predictive performance of the stationary distribution respectively.

5.2 Optimal compromise between Objective Posterior conditional distributions in Gaussian Process regression

Issues raised by objective prior elicitation for Gaussian processes

Let $Y(\mathbf{x})$, $\mathbf{x} \in \mathcal{D}$ be a real-valued random field on a bounded subset \mathcal{D} of \mathbb{R}^r . We assume Y is Gaussian with zero mean (or known mean) and with covariance of the form $\text{Cov}(Y(\mathbf{x}), Y(\mathbf{x}')) = \sigma^2 K_{\boldsymbol{\theta}}(\mathbf{x} - \mathbf{x}')$. σ^2 thus denotes the variance of the Gaussian Process and $\boldsymbol{\theta} \in (0, +\infty)^r$, hereafter named the “vector of correlation lengths”, is the vector of scaling parameters used by the chosen class of correlation kernels $K_{\boldsymbol{\theta}}$.

Consider a set of $n \in \mathbb{N}$ points $(\mathbf{x}^{(i)})_{i \in [1, n]}$ belonging to the domain \mathcal{D} . This set is called the design set and Y is observed at all points of this set. Let \mathbf{Y} be the Gaussian vector $(Y(\mathbf{x}^{(i)}))_{i \in [1, n]}$ and let $\boldsymbol{\Sigma}_{\boldsymbol{\theta}}$ be its correlation matrix: the distribution of \mathbf{Y} is therefore $\mathcal{N}(\mathbf{0}_n, \sigma^2 \boldsymbol{\Sigma}_{\boldsymbol{\theta}})$.

Let \mathbf{y} be the vector of observations. When applied to a matrix, $|\cdot|$ refers to its determinant.

With these notations, the likelihood of the parameters σ^2 and $\boldsymbol{\theta}$ is

$$L(\mathbf{y} \mid \sigma^2, \boldsymbol{\theta}) = \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{n}{2}} |\boldsymbol{\Sigma}_{\boldsymbol{\theta}}|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \mathbf{y}^\top \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} \mathbf{y} \right\}. \quad (5.1)$$

Conditional on $\boldsymbol{\theta}$, this is a scale model. The reference prior with parameter ordering $\sigma^2 \prec \boldsymbol{\theta}$ is therefore given by:

$$\pi(\sigma^2, \boldsymbol{\theta}) = \pi(\sigma^2 \mid \boldsymbol{\theta}) \pi(\boldsymbol{\theta}) \quad \text{with} \quad \pi(\sigma^2 \mid \boldsymbol{\theta}) \propto 1/\sigma^2. \quad (5.2)$$

The distribution $\pi(\sigma^2 \mid \boldsymbol{\theta})$ has infinite mass: it is an improper prior.

Let us integrate σ^2 out of the likelihood (5.1):

$$\begin{aligned} L^1(\mathbf{y} \mid \boldsymbol{\theta}) &\propto \int_0^\infty L(\mathbf{y} \mid \sigma^2, \boldsymbol{\theta}) \pi(\sigma^2 \mid \boldsymbol{\theta}) d(\sigma^2) \\ &\propto \int_0^\infty L(\mathbf{y} \mid \sigma^2, \boldsymbol{\theta}) (\sigma^2)^{-1} d(\sigma^2) = \left(\frac{2\pi^{\frac{n}{2}}}{\Gamma(\frac{n}{2})} \right)^{-1} |\boldsymbol{\Sigma}_{\boldsymbol{\theta}}|^{-\frac{1}{2}} (\mathbf{y}^\top \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} \mathbf{y})^{-\frac{n}{2}}. \end{aligned} \quad (5.3)$$

This model is the one from the remark under Example 3 in Chapter 2. The “marginal” reference prior on $\boldsymbol{\theta}$ is proportional to the square root of the determinant of the $r \times r$ matrix with (i, j) -th entry

$$\text{Tr} \left[\left(\frac{\partial}{\partial \theta_i} (\boldsymbol{\Sigma}_{\boldsymbol{\theta}}) \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} \right) \left(\frac{\partial}{\partial \theta_j} (\boldsymbol{\Sigma}_{\boldsymbol{\theta}}) \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} \right) \right] - \frac{1}{n} \text{Tr} \left[\frac{\partial}{\partial \theta_i} (\boldsymbol{\Sigma}_{\boldsymbol{\theta}}) \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} \right] \text{Tr} \left[\frac{\partial}{\partial \theta_j} (\boldsymbol{\Sigma}_{\boldsymbol{\theta}}) \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} \right]. \quad (5.4)$$

See also Ren et al. [2012] for an other way of reaching this conclusion. Their solution is based on astute algebraic considerations.

Viewing $\boldsymbol{\theta}$ as a single parameter has the disadvantage of requiring the use of a multidimensional Jeffreys-rule prior distribution, which may show the sort of undesirable behavior mentioned in the introduction, of which Example 4 from Chapter 2 is an illustration.

Alternatively, we could draw inspiration from the one-dimensional case in the following way. Suppose that we knew every entry of $\boldsymbol{\theta}$ except one, θ_i . Then, according to Equation (5.4), the prior density on θ_i knowing all entries θ_j ($j \neq i$) would be

$$\pi_i(\theta_i \mid \theta_j \forall j \neq i) \propto \sqrt{\text{Tr} \left[\left(\frac{\partial}{\partial \theta_i} (\boldsymbol{\Sigma}_\theta) \boldsymbol{\Sigma}_\theta^{-1} \right)^2 \right] - \frac{1}{n} \text{Tr} \left[\frac{\partial}{\partial \theta_i} (\boldsymbol{\Sigma}_\theta) \boldsymbol{\Sigma}_\theta^{-1} \right]^2}. \quad (5.5)$$

The density functions $\pi_i(\theta_i \mid \theta_j \forall j \neq i)$ ($i \in \llbracket 1, r \rrbracket$) define Markov kernels. Indeed, they are continuous with respect to the θ_j ($j \neq i$) and are probability densities with respect to the Lebesgue measure. They are unfortunately likely to violate the necessary condition for compatibility given by Equation (4.1).

Let us now consider the corresponding posterior conditional densities $\pi_i(\theta_i \mid \mathbf{y}, \theta_j \forall j \neq i)$ ($i \in \llbracket 1, r \rrbracket$). Just like their prior counterparts, they are likely to violate the necessary condition for compatibility. However, each of them represents our opinion about one parameter if all others were known. This is a setting where the results of Section 4.2 can be applied in order to find the optimal compromise between the Markov kernels $\mathbb{R}^{r-1} \times \mathcal{B}(\mathbb{R})$ they define. This optimal compromise will then be taken as posterior probability of the vector $\boldsymbol{\theta}$. In the following, we describe settings in which there exists a single Gibbs compromise between these Markov kernels. Theorem 4.9 then asserts it is the optimal compromise. We call this compromise the Gibbs reference posterior distribution because of its link to the reference posterior distribution in settings with a one-dimensional parameter $\boldsymbol{\theta}$.

However, even though we call it a “posterior” distribution, it is unclear whether there exists a prior distribution from which it could be derived using Bayes’ rule. Denote by $\pi_G(\boldsymbol{\theta} \mid \mathbf{y})$ the probability density with respect to the Lebesgue measure of the Gibbs reference posterior distribution. Bayes’ rule requires that in case a (proper or improper) prior density $\pi_G(\boldsymbol{\theta})$ exists, there should also exist a function $\tilde{L}(\mathbf{y})$ such that, for almost every $\boldsymbol{\theta} \in \mathbb{R}^r$ in the sense of the Lebesgue measure,

$$\frac{\pi_G(\boldsymbol{\theta} \mid \mathbf{y})}{L^1(\mathbf{y} \mid \boldsymbol{\theta})} = \frac{\pi_G(\boldsymbol{\theta})}{\tilde{L}(\mathbf{y})}. \quad (5.6)$$

As we have no explicit expression of $\pi_G(\boldsymbol{\theta} \mid \mathbf{y})$, we have no way to check whether Equation (5.6) holds or not.

In this section, we establish that, whenever Matérn anisotropic geometric or tensorized kernels with known smoothness parameter ν are used, under certain conditions to be detailed later, there exists a unique Gibbs compromise between the reference posterior conditionals, which thanks to Theorem 4.9 is the optimal compromise. Henceforth, it will be called “Gibbs reference posterior distribution”, even though this “posterior” has not been derived from a prior distribution using the Bayes rule.

All proofs for this section can be found in Appendix 5.A.

Definitions

In this chapter, we use the following convention for the Fourier transform: the Fourier transform \widehat{g} of a smooth function $g : \mathbb{R}^r \rightarrow \mathbb{R}$ verifies $g(\mathbf{x}) = \int_{\mathbb{R}^r} \widehat{g}(\boldsymbol{\omega}) e^{i\langle \boldsymbol{\omega} | \mathbf{x} \rangle} d\boldsymbol{\omega}$ and $\widehat{g}(\boldsymbol{\omega}) = (2\pi)^{-r} \int_{\mathbb{R}^r} g(\mathbf{x}) e^{-i\langle \boldsymbol{\omega} | \mathbf{x} \rangle} d\mathbf{x}$.

Let us set up a few notations.

- (a) \mathcal{K}_ν is the modified Bessel function of second kind with parameter ν ;
 (b) $K_{r,\nu}$ is the r -dimensional Matérn isotropic covariance kernel with variance 1, correlation length 1 and smoothness $\nu \in (0, +\infty)$ and $\widehat{K}_{r,\nu}$ is its Fourier transform:

$$(i) \quad \forall \mathbf{x} \in \mathbb{R}^r, \quad K_{r,\nu}(\mathbf{x}) = \frac{1}{\Gamma(\nu)2^{\nu-1}} (2\sqrt{\nu}\|\mathbf{x}\|)^\nu \mathcal{K}_\nu(2\sqrt{\nu}\|\mathbf{x}\|) ; \quad (5.7)$$

$$(ii) \quad \forall \boldsymbol{\omega} \in \mathbb{R}^r, \quad \widehat{K}_{r,\nu}(\boldsymbol{\omega}) = \frac{M_r(\nu)}{(\|\boldsymbol{\omega}\|^2 + 4\nu)^{\nu + \frac{r}{2}}} \text{ with } M_r(\nu) = \frac{\Gamma(\nu + \frac{r}{2})(2\sqrt{\nu})^{2\nu}}{\pi^{\frac{r}{2}}\Gamma(\nu)}. \quad (5.8)$$

- (c) $K_{r,\nu}^{tens}$ is the r -dimensional Matérn tensorized covariance kernel with variance 1, correlation length 1 and smoothness $\nu \in \mathbb{R}_+$ and $\widehat{K}_{r,\nu}^{tens}$ is its Fourier transform:

$$(i) \quad \forall \mathbf{x} \in \mathbb{R}^r, \quad K_{r,\nu}^{tens}(\mathbf{x}) = \prod_{j=1}^r K_{1,\nu}(\mathbf{x}_j) ; \quad (5.9)$$

$$(ii) \quad \forall \boldsymbol{\omega} \in \mathbb{R}^r, \quad \widehat{K}_{r,\nu}^{tens}(\boldsymbol{\omega}) = \prod_{j=1}^r \widehat{K}_{1,\nu}(\boldsymbol{\omega}_j). \quad (5.10)$$

- (d) if $\mathbf{t} \in \mathbb{R}^r$, $\frac{\mathbf{t}}{\boldsymbol{\theta}} = \left(\frac{t_1}{\theta_1}, \dots, \frac{t_r}{\theta_r}\right)$ and $\mathbf{t}\boldsymbol{\mu} = (t_1\mu_1, \dots, t_r\mu_r)$.

We define the Matérn geometric anisotropic covariance kernel with variance parameter σ^2 , correlation lengths $\boldsymbol{\theta}$ (resp. inverse correlation lengths $\boldsymbol{\mu}$) and smoothness ν as the function $\mathbf{x} \mapsto \sigma^2 K_{r,\nu}\left(\frac{\mathbf{x}}{\boldsymbol{\theta}}\right)$ (resp. $\mathbf{x} \mapsto \sigma^2 K_{r,\nu}(\mathbf{x}\boldsymbol{\mu})$).

Similarly, we define the Matérn tensorized covariance kernel with variance parameter σ^2 , correlation lengths $\boldsymbol{\theta}$ (resp. inverse correlation lengths $\boldsymbol{\mu}$) and smoothness ν as the function $\mathbf{x} \mapsto \sigma^2 K_{r,\nu}^{tens}\left(\frac{\mathbf{x}}{\boldsymbol{\theta}}\right)$ (resp. $\mathbf{x} \mapsto \sigma^2 K_{r,\nu}^{tens}(\mathbf{x}\boldsymbol{\mu})$).

Thanks to Proposition 4.10, we may choose any parametrization we wish for the Matérn correlation kernels. We have found that the parametrization involving inverse correlation lengths makes proofs easier.

Several key passages in the proofs (to be found in Appendix 5.A) involve a technical assumption on the design set:

Definition 5.1. *A design set is said to have coordinate-distinct points, or simply to be coordinate-distinct, if for any distinct points in the set \mathbf{x} and \mathbf{x}' , every component of the vector $\mathbf{x} - \mathbf{x}'$ differs from 0.*

Randomly sampled design sets almost surely have coordinate-distinct points – for instance Latin Hypercube Sampling. Cartesian product design sets, however, do not.

Main result

The result is valid for Simple Kriging models with the following characteristics:

- (a) the design set contains n coordinate-distinct points in \mathbb{R}^r (n and r are positive integers);
- (b) the covariance function is Matérn anisotropic geometric or tensorized with variance parameter $\sigma^2 > 0$, smoothness parameter ν and vector of correlation lengths (resp. inverse correlation lengths) $\boldsymbol{\theta} \in (0, +\infty)^r$ (resp. $\boldsymbol{\mu} \in (0, +\infty)^r$);
- (c) one of the following conditions is verified:
 - (i) $\nu \in (0, 1)$ and $n > 1$ and the Matérn kernel is tensorized;
 - (ii) $\nu \in (1, 2)$ and $n > r + 2$;
 - (iii) $\nu \in (2, 3)$ and $n > r(r + 1)/2 + 2r + 3$.

Theorem 5.2. *In a Simple Kriging model with the characteristics described above, there exists a hyperplane \mathcal{H} of \mathbb{R}^n such that, for any $\mathbf{y} \in \mathbb{R}^n \setminus \mathcal{H}$, there exists a unique Gibbs compromise $\pi_G(\boldsymbol{\theta}|\mathbf{y})$ (resp. $\pi_G(\boldsymbol{\mu}|\mathbf{y})$) between the reference posterior conditionals $\pi_i(\theta_i|\mathbf{y}, \boldsymbol{\theta}_{-i})$ (resp. $\pi_i(\mu_i|\mathbf{y}, \boldsymbol{\mu}_{-i})$). It is the unique stationary distribution of the Markov kernel $P_{\mathbf{y}} : (0, +\infty)^r \times \mathcal{B}((0, +\infty)^r) \rightarrow [0, 1]$ defined by*

$$P_{\mathbf{y}}(\boldsymbol{\theta}^{(0)}, d\boldsymbol{\theta}) = \frac{1}{r} \sum_{i=1}^r \pi_i(\theta_i|\mathbf{y}, \boldsymbol{\theta}_{-i}^{(0)}) d\theta_i \delta_{\boldsymbol{\theta}_{-i}^{(0)}}(d\boldsymbol{\theta}_{-i})$$

$$(resp. P_{\mathbf{y}}(\boldsymbol{\mu}^{(0)}, d\boldsymbol{\mu}) = \frac{1}{r} \sum_{i=1}^r \pi_i(\mu_i|\mathbf{y}, \boldsymbol{\mu}_{-i}^{(0)}) d\mu_i \delta_{\boldsymbol{\mu}_{-i}^{(0)}}(d\boldsymbol{\mu}_{-i})).$$

The Markov kernel $P_{\mathbf{y}}$ is uniformly ergodic. This means that, denoting by $\|\cdot\|_{TV}$ the total variation norm,

$$\lim_{n \rightarrow \infty} \sup_{\boldsymbol{\theta}^{(0)} \in (0, +\infty)^r} \|P_{\mathbf{y}}^n(\boldsymbol{\theta}^{(0)}, \cdot) - \pi_G(\cdot|\mathbf{y})\|_{TV} = 0$$

$$(resp. \lim_{n \rightarrow \infty} \sup_{\boldsymbol{\mu}^{(0)} \in (0, +\infty)^r} \|P_{\mathbf{y}}^n(\boldsymbol{\mu}^{(0)}, \cdot) - \pi_G(\cdot|\mathbf{y})\|_{TV} = 0). \quad (5.11)$$

Remark. The reference posterior conditionals are invariant by reparametrization, so the Markov kernel $P_{\mathbf{y}}$ does not depend on whether the chosen parametrization uses correlation lengths $\boldsymbol{\theta}$ or inverse correlation lengths $\boldsymbol{\mu}$. Due to Proposition 4.10, the Gibbs compromise does not either. The parametrization using inverse correlation lengths $\boldsymbol{\mu}$ is more convenient for proving this theorem, however, so we use it exclusively in the rest of this section.

Notice that in such a Kriging model, the vector of observations \mathbf{y} almost surely belongs to $\mathbb{R}^n \setminus \mathcal{H}$, so this assumption is of no practical consequence. Theorem 5.2 therefore asserts that the Gibbs compromise between the incompatible conditionals $\pi_i(\mu_i|\mathbf{y}, \boldsymbol{\mu}_{-i})$ exists, is unique, and can be sampled from using Potentially Incompatible Gibbs Sampling (PIGS). In the following, it is called ‘‘Gibbs reference posterior distribution’’.

Using the Gibbs reference posterior distribution

Let \mathbf{x}_0 be a point in the domain \mathcal{D} that does not belong to the design set. Denote by $\boldsymbol{\Sigma}_{\boldsymbol{\theta}, 0}$, the correlation matrix between $Y(\mathbf{x}_0)$ and \mathbf{Y} , and by $\boldsymbol{\Sigma}_{\boldsymbol{\theta}, \cdot, 0}$ its transpose the correlation matrix between \mathbf{Y} and $Y(\mathbf{x}_0)$.

Theorem 4.1.2. (case 4) of Santner et al. [2003] provides this useful result for prediction:

Proposition 5.3. *Conditionally to $\mathbf{Y} = \mathbf{y}$ and assuming $\boldsymbol{\theta}$ is known, the random variable Z_0 defined below follows the Student t -distribution with n degrees of freedom.*

$$Z_0 := \sqrt{\frac{n}{\mathbf{y}^\top \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} \mathbf{y}}} \frac{Y(\mathbf{x}_0) - \boldsymbol{\Sigma}_{\boldsymbol{\theta},0} \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} \mathbf{y}}{\sqrt{1 - \boldsymbol{\Sigma}_{\boldsymbol{\theta},0} \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta},0}}}$$

Remark. If n exceeds 30, it is usually accepted that the Student t -distribution with n degrees of freedom can be approximated by the standard Normal distribution. As this threshold should be exceeded in practical cases, we would recommend performing all computations as though the Student t -distribution were Normal.

In practice, the distribution of $Y(\mathbf{x}_0)$ conditionally to $\mathbf{Y} = \mathbf{y}$ when $\boldsymbol{\theta}$ is unknown can be obtained once from the Gibbs reference posterior distribution on $\boldsymbol{\theta}$ has been sampled. Its cdf can be approximated by averaging the cdfs of the Student t -distributions (or their Normal approximations) corresponding to every point in the sample.

5.3 Comparisons between the MLE and MAP estimators

To illustrate the inferential performance of the Gibbs reference posterior distribution, let us introduce the Maximum A Posteriori estimator (MAP). It takes the value of $\boldsymbol{\theta}$ where the density with respect to the Lebesgue measure of the Gibbs reference posterior distribution is largest. We contrast it with the Maximum Likelihood Estimator (MLE) which does the same with the likelihood function.

Methodology

In this section, we compare the MLE and MAP estimators for accuracy and robustness.

Our test cases are 3-dimensional Gaussian Processes with Matérn anisotropic geometric correlation kernels with smoothness $5/2$. Their mean is the null function, which only leaves us with the matter of estimating their correlation length for each dimension.

We use uniform designs: our observation points are randomly generated according to the uniform distribution on a cube with side length 1.

In order to measure the performance of our estimators, we define a suitable distance between two vectors of correlation lengths. Then the error of an estimator is defined as its distance to the “true” vector of correlation lengths.

Let g be the function such that for any t in $(-1, 1)$, $g(t) = \operatorname{arctanh}(t)$ and $g(-1) = g(1) = 0$. We use the convention that, for any matrix \mathbf{M} with elements in $[0, 1]$, $g(\mathbf{M})$ is the matrix resulting from applying g to every element of \mathbf{M} .

Definition 5.4. *For a given design set, the distance between two vectors of correlation lengths $\boldsymbol{\theta}^1$ and $\boldsymbol{\theta}^2$ is $\|g(\boldsymbol{\Sigma}_{\boldsymbol{\theta}^1}) - g(\boldsymbol{\Sigma}_{\boldsymbol{\theta}^2})\|$, where $\|\cdot\|$ denotes the Frobenius norm.*

This distance involves applying the Fisher transformation [Hotelling, 1953] (that is, the inverse hyperbolic tangent function) to every (non-1) correlation coefficient in both associated correlation matrices. This is a variance-stabilizing transformation. For any random variables U and V following the normal distribution with mean 0 and variance 1, let ρ denote the correlation coefficient between U and V ($-1 < \rho < 1$). If (U_i, V_i) ($1 \leq i \leq N$) are independent copies of (U, V) , then $\hat{\rho} = \sum_{i=1}^N U_i V_i / n$ is a random variable and $\operatorname{arctanh}(\hat{\rho})$ follows

the normal distribution with mean $\operatorname{arctanh}(\rho)$ and variance $1/(N - 3)$. So the variance of $\operatorname{arctanh}(\hat{\rho})$ does not depend on ρ , whereas the variance of $\hat{\rho}$ does and goes to zero for $|\rho| \rightarrow 1$. Involving the Fisher transformation in the definition of the distance between two vectors of correlation lengths is therefore a way to assert that vectors of correlation lengths can be far apart even if they both lead to highly correlated observations.

This allows us to make sure errors made when estimating near-1 correlation coefficients are no less taken into account than errors made when estimating near-0 correlation coefficients.

Let us choose a “true” vector of correlation lengths (and also a variance parameter, but this parameter has no effect on either the MLE or the MAP). Then we need to:

1. Sample n points randomly according to the uniform distribution on the unit cube (in the following, $n = 30$).
2. Generate the observations of the Gaussian Process at the sampled points according to the selected “true” variance and correlation lengths.
3. Sample the vector of correlation lengths according to the Gibbs reference posterior distribution $\pi_G(\boldsymbol{\theta}|\mathbf{y})$ through PIGS.
4. Compute the MLE and the MAP of the vector of correlation lengths and their errors.
5. Repeat steps 1 to 4 $m - 1$ times (in the following, $m = 500$).

This method allows us to derive an approximate distribution of the errors of both estimators when both the realization of the Gaussian Process and the design set vary. Thus we get to test the robustness of both estimators versus the variability of both the Gaussian Process and the choice of design set.

Results

This subsection provides results obtained on 3-dimensional Gaussian Processes with null mean function and Matérn anisotropic geometric correlation kernels with smoothness $5/2$. The results are divided by “true” vectors of correlation lengths. In each case, we give in Table 5.1 the empirical Root Mean Square Errors (RMSEs) of both MLE and MAP estimators as functions of varying instances of the Gaussian Process and uniform design sets on the unit cube.

Most of the “true” vectors of correlation lengths featured in Table 5.1 were selected in a way to showcase the behavior of both estimators in strongly anisotropic cases, but one (0.5 - 0.5 - 0.5) also showcases their behavior if the true kernel is actually isotropic. And the final one (0.8 - 1 - 0.9) is used to illustrate the performance in the case of a strongly correlated Gaussian Process: this case is fundamentally different from all others, because the Matérn anisotropic geometric family of correlation kernels is designed in such a way that the correlation length with greatest influence is the lowest. Informally speaking, it is enough for one correlation length to be near zero to make the whole process very uncorrelated, even should all other correlation lengths be very high.

In all studied cases, the MAP estimator was more robust than the MLE estimator: its RMSE (Root Mean Square Error) was between 9 and 15% lower, as showcased in Table 5.1.

Corr. lengths	MLE	MAP	- (%)
0.4 – 0.8 – 0.2	3.49	2.97	15
0.5 – 0.5 – 0.5	4.00	3.46	13
0.7 – 1.3 – 0.4	4.02	3.64	9
0.8 – 0.3 – 0.6	3.75	3.26	13
0.8 – 1.0 – 0.9	4.65	4.18	10

Table 5.1 – RMSE (where the error is measured in terms of the distance in Definition 5.4) of the MLE and MAP estimators for several “true” vectors of correlation lengths. The last column displays in percents the decrease of the RMSE of the MAP estimator with respect to the MLE.

To get a better sense of the distribution of the error when the design set and the realization of the Gaussian Process vary, we give in Figure 5.1 violin plots of the errors in the two most extreme case: very low correlation (0.4 – 0.8 – 0.2) and very high correlation (0.8 – 1.0 – 0.9)

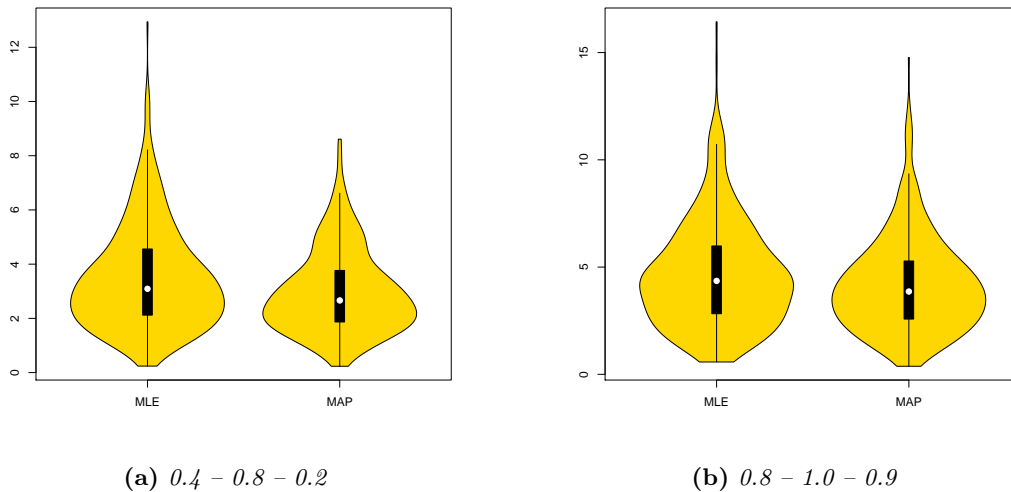


Figure 5.1 – Violin plots of the error of the MLE and MAP estimators with respect to a design set following the uniform distribution and a Gaussian Process with correlation lengths 0.4 – 0.8 – 0.2 (left) and 0.8 – 1.0 – 0.9 (right).

5.4 Comparison of the predictive distributions associated with the estimators (MLE and MAP) and the full posterior distribution

Methodology

We use the same test cases as before. In this section, our goal is to assess the accuracy of prediction intervals associated with both estimators and with the full posterior distribution. We consider 95% intervals: the lower bound is the 2.5% quantile and the upper bound the 97.5% quantile of predictive distribution $\hat{P}_{MLE}(\mathbf{y}_0 | \mathbf{y})$, $\hat{P}_{MAP}(\mathbf{y}_0 | \mathbf{y})$ and $P(\mathbf{y}_0 | \mathbf{y})$. For the sake of comprehensiveness, we also consider predictive intervals associated with the “true” predictive distribution $L(\mathbf{y}_0 | \mathbf{y}, \sigma^2, \boldsymbol{\theta})$, which is the predictive distribution we would use if we knew the correct values of the parameters σ^2 and $\boldsymbol{\theta}$.

Let us choose a “true” vector of correlation lengths $\boldsymbol{\theta}$ (and also a variance parameter σ^2 , but this parameter has no effect on predictive accuracy). Then we do the following:

1. Sample n observation points randomly according to the uniform distribution on the unit cube (in the following, $n = 30$).
2. Generate the observations of the Gaussian Process at the sampled points according to the selected “true” variance and correlation lengths.
3. Sample the vector of correlation lengths according to the Gibbs reference posterior distribution $\pi_G(\boldsymbol{\theta} | \mathbf{y})$ through PIGS.
4. Compute the MLE and the MAP of the vector of correlation lengths.
5. Sample n_0 test points randomly according to the uniform distribution on the unit cube (in the following, $n_0 = 100$).
6. At each point, determine the 95% prediction intervals derived from $L(\mathbf{y}_0 | \mathbf{y}, \sigma^2, \boldsymbol{\theta})$ (σ^2 and $\boldsymbol{\theta}$ being the “true” parameters), $\hat{P}_{MLE}(\mathbf{y}_0 | \mathbf{y})$, $\hat{P}_{MAP}(\mathbf{y}_0 | \mathbf{y})$ and $P(\mathbf{y}_0 | \mathbf{y})$.
7. Generate the values of the Gaussian Process at the newly sampled points (naturally, do this conditionally to the previously generated observations).
8. Count the number of points within the prediction intervals derived of each of the four predictive distributions. Divide the counts by n_0 : this yields four *coverages* corresponding to each type of predictive intervals. Also compute the *mean length* of every type of prediction interval.
9. Repeat steps 1 to 8 $m - 1$ times (in the following, $m = 500$).

Results

There is no reason for individual coverages of 95% predictive intervals given by the predictive distribution to be equal to 95%. Recall that any coverage is given for a unique realization of the Gaussian Process, and that the values of this process at different points are correlated. If the predictive interval at some point fails to cover the true value at this point, it is likely that predictive intervals at neighboring points will also fail to cover the true values at those points, even though the nominal value is 95% everywhere. Conversely, if it actually covers the true value, then prediction intervals at neighboring points are more than 95% likely to cover their true values.

In short, prediction intervals give information that is only valid if understood to refer to what can be guessed on the sole basis of the observations made at the design points, which is why coverages for individual realizations of the Gaussian Process are not necessarily 95% *even if the predictive distribution is perfectly accurate* (*i.e.* based on the true values of σ^2 and θ).

However, *if the predictive distribution is perfectly accurate*, then the average of the coverages is the nominal value: 95%. It is thus interesting to compute the average of the coverages for all predictive distributions, whether they are based on the MLE or MAP estimator, or on the full posterior distribution (hereafter noted FPD). In the above described methodology, the average was taken over the realizations of the Gaussian Process with the chosen true parameters and over all design sets with n design points. The results below are obtained in this way.

The results given in Table 5.2 show that using the full posterior distribution (FPD) to derive the predictive distribution is the best possible choice from a frequentist point of view as the nominal value is nearly matched by the average coverage. Predictive intervals derived from the MAP estimator do not perform as well, and predictive intervals derived from the MLE perform even worse.

Corr. lengths	True	MLE	MAP	FPD
0.4 – 0.8 – 0.2	0.95	0.88	0.91	0.95
0.5 – 0.5 – 0.5	0.95	0.89	0.90	0.94
0.7 – 1.3 – 0.4	0.95	0.90	0.92	0.95
0.8 – 0.3 – 0.6	0.95	0.89	0.91	0.95
0.8 – 1.0 – 0.9	0.95	0.90	0.92	0.94

Table 5.2 – Average with respect to randomly sampled design sets and realizations of the Gaussian Process (with variance parameter 1 and smoothness parameter 5/2) of the coverage of 95% Prediction Intervals across the sample space. “True” stands for the prediction based on the knowledge of the true variance parameter and the true vector of correlation lengths.

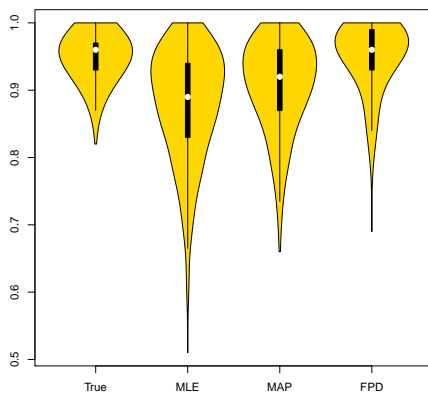
Let us now focus on the average (with respect to the uniform design sets and realizations of the Gaussian Process) of the mean (over the test set for a given realization of the Gaussian Process and a given uniform design set) length of prediction intervals. The results are given in Table 5.3, where the figures between parentheses give the increase or decrease (in percents) of the average mean length when compared to the average mean length of prediction intervals obtained using the true values of the parameters.

Corr. lengths	True	MLE	MAP	FPD
0.4 – 0.8 – 0.2	2.23	2.05 (-8)	2.13 (-4)	2.59 (+16)
0.5 – 0.5 – 0.5	1.69	1.55 (-8)	1.58 (-6)	1.84 (+9)
0.7 – 1.3 – 0.4	1.09	1.02 (-6)	1.07 (-2)	1.21 (+11)
0.8 – 0.3 – 0.6	1.63	1.51 (-7)	1.56 (-4)	1.82 (+12)
0.8 – 1.0 – 0.9	0.71	0.66 (-7)	0.69 (-3)	0.76 (+8)

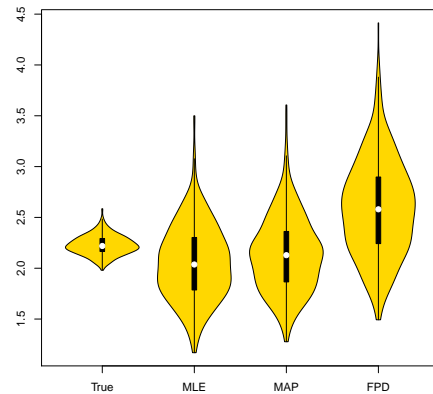
Table 5.3 – Average with respect to randomly sampled design sets and realizations of the Gaussian Process (with variance parameter 1 and smoothness parameter 5/2) of the mean length of 95% Prediction Intervals across the sample space. The numbers in parentheses represent in percents the increase when using the MLE/MAP/FPD instead of the “true” vector of correlation lengths and variance parameter.

Predictive intervals derived from the full posterior distribution (FPD) are on average the largest, but not much larger than predictive intervals derived using the true parameters. In the tests we conducted, they seemed on average to be larger by about one fifth at worst. Predictive intervals derived from the MLE and MAP estimators are on average shorter than those derived from the true parameters. This can be interpreted as an under-estimation of the uncertainty of the prediction when fixing the vector of correlation lengths to the most likely value given the observations, and this can explain the low observed coverage in Table 5.2.

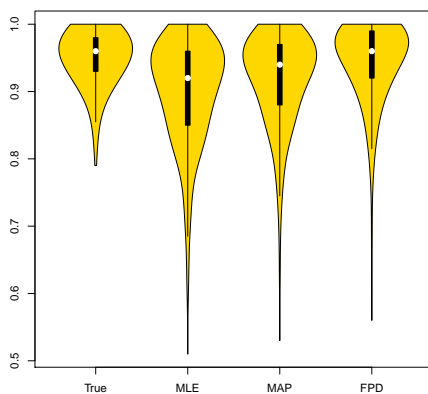
In Figure 5.2, we give violin plots of coverage and mean length of Prediction Intervals in the two most extreme cases: correlation lengths $0.4 - 0.8 - 0.2$ (very low correlation) and $0.8 - 1.0 - 0.9$ (very high correlation). The results are similar and illustrate the fact that the FPD gives larger intervals in order to reach the derived coverage value.



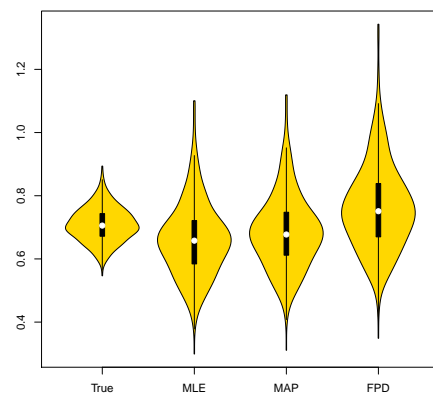
(a) Coverage for $0.4 - 0.8 - 0.2$



(b) Mean length for $0.4 - 0.8 - 0.2$



(c) Coverage for $0.8 - 1.0 - 0.9$



(d) Mean length for $0.8 - 1.0 - 0.9$

Figure 5.2 – Violin plots of the coverage (left) and mean length (right) of Prediction Intervals with respect to a design set following the uniform distribution and a Gaussian Process with correlation lengths $0.4 - 0.8 - 0.2$ (top) and $0.8 - 1.0 - 0.9$ (bottom).

A higher-dimensional case

In this subsection, we emulate using Simple Kriging the 10-dimensional Ackley function:

$$A(\mathbf{x}) = 20 + \exp(1) - 20 \exp \left(-0.2 \sqrt{\frac{1}{10} \sum_{i=1}^{10} x_i^2} \right) - \exp \left(\frac{1}{10} \sum_{i=1}^{10} \cos(2\pi x_i) \right). \quad (5.12)$$

The goal in this section is to emulate the Ackley function on the unit hypercube $[0, 1]^{10}$ using design sets with 100 observation points. Although the impact of the design set type is not the focus of this study, we present the results with a randomly chosen design according to the Uniform distribution on the domain $[0, 1]^{10}$, a design obtained through Latin Hypercube Sampling (LHS), and a design obtained through LHS and subsequently optimized to maximize the minimum distance between two points. The Simple Kriging model uses the null function as mean function and the Matérn anisotropic geometric covariance kernel family with smoothness parameter $5/2$. The Gibbs reference posterior distribution is accessed through a sample of 1000 points. The conditional densities are sampled using the Metropolis algorithm with normal instrumental density with standard deviation 0.4 and a 100-step burn-in period.

To evaluate the performance of prediction intervals, we follow steps 3, 4, 5, 6 and 8 of the method presented in this section (step 7 is skipped as the “values of the Gaussian process” are naturally the values of the Ackley function) with $n_0 = 1000$. The results are presented in Tables 5.4 and 5.5.

Design set type	MLE	MAP	FPD
Unoptimized LHS	0.89	0.92	0.93
Optimized LHS	0.74	0.76	0.80
Random design	0.87	0.88	0.91

Table 5.4 – Coverage of 95 % prediction intervals when emulating the Ackley function on the unit hypercube using a Gaussian Process with null mean function and a Matérn anisotropic geometric covariance kernel with smoothness $5/2$, unknown variance parameter and unknown vector of correlation lengths. The design sets contain 100 points.

As is shown in Table 5.4, prediction intervals derived using the Full Posterior Distribution perform better than those derived from the MAP, which themselves perform better than those derived from the MLE. This order of performance is the same regardless of the type of design set, although the optimized design set leads to much worse performances on average for prediction intervals than unoptimized designs. The latter fact is not surprising since space-filling designs ensure that no two points can be very close to each other, which makes it harder to determine the correlation lengths.

Design set type	MLE	MAP	FPD
Unoptimized LHS	0.31	0.33	0.35
Optimized LHS	0.24	0.24	0.28
Random design	0.28	0.29	0.32

Table 5.5 – Mean length of 95 % prediction intervals when emulating the Ackley function on the unit hypercube using a Gaussian Process with null mean function and a Matérn anisotropic geometric covariance kernel with smoothness $5/2$, unknown variance parameter and unknown vector of correlation lengths. The design sets contain 100 points.

As expected, prediction intervals derived from the Full Posterior Distribution are on average longer than those derived from the MAP and a fortiori the MLE. Notice that prediction intervals are on average shorter with the optimized design set, which explains the poorer performances in terms of coverage.

5.5 Conclusion and Perspectives

We provided theoretical foundation to the claim that the stationary distribution of the Markov chain underlying PIGS with random scanning order is the optimal compromise between the potentially incompatible conditional distributions.

Although further investigation is needed to fully understand the properties of the optimal compromise, its invariance by reparametrization and its respect of pairwise independence show that it preserves important features of the conditional distributions.

This construction suggests a framework for deriving a new objective posterior distribution based on the conditionals yielded by the reference prior theory on Simple Kriging parameters. Applying this framework to Matérn anisotropic kernels, we showed prediction to have good frequentist properties.

The next step, which is taken in Chapter 6, is to extend this framework to Universal Kriging, where instead of being known, the mean function is only assumed to be a linear combination of known functions f_1, \dots, f_p . The linear coefficients β_1, \dots, β_p are then considered parameters of the model. This extension is of practical relevance, because the mean function can rarely be considered known. It can be done in the same way Berger et al. [2001] extended the reference prior from the Simple Kriging to the Universal Kriging framework: they used the flat improper prior as joint prior on β_1, \dots, β_p conditional to σ^2 and $\boldsymbol{\theta}$ and used it to integrate β_1, \dots, β_p out of the likelihood function, and then proceeded to derive the reference prior on σ^2 and $\boldsymbol{\theta}$ with respect to the integrated likelihood.

A further extension would involve deriving an objective prior on the smoothness parameter ν . In this endeavor, one should take into account the relationship between correlation length $\boldsymbol{\theta}$ and smoothness ν . Unfortunately, asymptotic theory is not of much help in this regard, as Anderes [2010] shows that provided the spatial domain \mathcal{D} is of dimension at least 5, then all parameters of the Matérn anisotropic geometric kernel are microergodic (Zhang [2004] shows this to be untrue for spatial domains of dimension 1, 2 or 3, but the non-microergodic parameters are σ^2 and $\boldsymbol{\theta}$, not ν). This means that the Gaussian measures on \mathcal{D} corresponding to Gaussian Processes with two different smoothness parameters are orthogonal, which suggests that there exists a consistent estimator (the MLE possibly). Stein [1999] (section 6.6) considers the Fisher information on $\boldsymbol{\theta}$ and ν , and gives examples (with a one-dimensional sample space \mathcal{D}) showing that the Fisher information on these parameters depends a lot on the design set. Fisher information relative to the smoothness parameter ν increases when design points are chosen to be close to one another (relative to the "true" correlation length $\boldsymbol{\theta}$), whereas Fisher information relative to correlation length $\boldsymbol{\theta}$ is maximized for design points that are farther apart. This, according to him, is coherent with the fact that $\boldsymbol{\theta}$ has greater influence on the low frequency behavior of the Matérn kernel while ν has greater influence on its high frequency behavior. This also suggests to us that the smoothness parameter ν , like the variance parameter σ^2 , can only be meaningfully estimated if the vector of correlation

lengths θ is known. Otherwise, the estimator could hardly tell which design points are close to each other, which intuitively seems a prerequisite to evaluating the smoothness of the process. If we wish to apply the reference prior algorithm to the case where ν is unknown, we should thus probably derive the reference prior on ν conditional to θ .

Appendix 5.A Proofs of Section 5.2

The following holds where there is no mention of the contrary. When applied to a vector, $\|\cdot\|$ denotes the Euclidean norm and when applied to a matrix, it denotes the Frobenius norm. The choice of norm does not matter much because in finite-dimensional vector spaces, all norms are equivalent.

Differentiating the Matérn correlation kernel

Lemma 5.5. *The partial derivative with respect to μ_i of the Matérn tensorized kernel of variance σ^2 , smoothness ν and inverse correlation length vector $\boldsymbol{\mu}$ is:*

$$\frac{\partial}{\partial \mu_i} (\sigma^2 K_{r,\nu}^{tens}(\mathbf{x}\boldsymbol{\mu})) = -\frac{\sigma^2 (2\sqrt{\nu})^2}{\Gamma(\nu) 2^{\nu-1}} |x_i|^2 \mu_i (2\sqrt{\nu} |x_i| \mu_i)^{\nu-1} \mathcal{K}_{\nu-1}(2\sqrt{\nu} |x_i| \mu_i) \prod_{j \neq i} K_{1,\nu}(|x_j| \mu_j). \quad (5.13)$$

This can be rewritten as:

$$\frac{\partial}{\partial \mu_i} (\sigma^2 K_{r,\nu}^{tens}(\mathbf{x}\boldsymbol{\mu})) = \begin{cases} \sigma^2 \frac{2\nu}{\nu-1} |x_i|^2 \mu_i K_{1,\nu-1}(|x_i| \mu_i) \prod_{j \neq i} K_{1,\nu}(|x_j| \mu_j) & \text{if } \nu > 1 \\ \sigma^2 4 |x_i|^2 \mu_i \mathcal{K}_0(2|x_i| \mu_i) \prod_{j \neq i} K_{1,\nu}(|x_j| \mu_j) & \text{if } \nu = 1 \\ \sigma^2 \frac{2\nu^\nu \Gamma(1-\nu)}{\Gamma(\nu)} |x_i|^{2\nu} \mu_i^{2\nu-1} K_{1,1-\nu}(|x_i| \mu_i) \prod_{j \neq i} K_{1,\nu}(|x_j| \mu_j) & \text{if } \nu < 1. \end{cases} \quad (5.14)$$

Proof. The first assertion is a simple matter of differentiating Equation (5.9). In the following calculation, the fourth line is given by formula 9.6.28 (page 376) in Abramowitz and Stegun [1964].

$$\begin{aligned} \frac{\partial}{\partial \mu_i} (\sigma^2 K_{r,\nu}^{tens}(\mathbf{x}\boldsymbol{\mu})) &= \sigma^2 \frac{\partial}{\partial \mu_i} (K_{1,\nu}(x_i \mu_i)) \prod_{j \neq i} K_{1,\nu}(|x_j| \mu_j) \\ &= \sigma^2 x_i (K'_{1,\nu}(x_i \mu_i)) \prod_{j \neq i} K_{1,\nu}(|x_j| \mu_j) \\ &= \sigma^2 x_i \left(\frac{2\sqrt{\nu}}{\Gamma(\nu) 2^{\nu-1}} \frac{d}{dy} \Big|_{y=2\sqrt{\nu} x_i \mu_i} [y^\nu \mathcal{K}_\nu(y)] \right) \prod_{j \neq i} K_{1,\nu}(|x_j| \mu_j) \\ &= \sigma^2 x_i \left(\frac{2\sqrt{\nu}}{\Gamma(\nu) 2^{\nu-1}} [-y \cdot y^{\nu-1} \mathcal{K}_{\nu-1}(y)]_{y=2\sqrt{\nu} x_i \mu_i} \right) \prod_{j \neq i} K_{1,\nu}(|x_j| \mu_j). \end{aligned} \quad (5.15)$$

From there, Equation (5.13) follows immediately. Rewriting it in the form given in (5.14) only requires us to recall $\Gamma(\nu) = (\nu-1)\Gamma(\nu-1)$ (case $\nu > 1$), $\Gamma(1) = 1$ (case $\nu = 1$) and $\mathcal{K}_{\nu-1} = \mathcal{K}_{1-\nu}$ (case $\nu < 1$). \square

Lemma 5.6. *The partial derivative with respect to μ_i of the Matérn geometric anisotropic kernel of variance σ^2 , smoothness ν and inverse correlation length vector $\boldsymbol{\mu}$ is:*

$$\frac{\partial}{\partial \mu_i} (\sigma^2 K_{r,\nu}(\mathbf{x}\boldsymbol{\mu})) = \frac{\sigma^2 (2\sqrt{\nu})^2}{\Gamma(\nu) 2^{\nu-1}} |x_i|^2 \mu_i (2\sqrt{\nu} \|\mathbf{x}\boldsymbol{\mu}\|)^{\nu-1} \mathcal{K}_{\nu-1}(2\sqrt{\nu} \|\mathbf{x}\boldsymbol{\mu}\|). \quad (5.16)$$

This can be rewritten as:

$$\frac{\partial}{\partial \mu_i} (\sigma^2 K_{r,\nu}(\mathbf{x}\boldsymbol{\mu})) = \begin{cases} \sigma^2 \frac{2\nu}{\nu-1} |x_i|^2 \mu_i K_{1,\nu-1}(\|\mathbf{x}\boldsymbol{\mu}\|) & \text{if } \nu > 1 \\ \sigma^2 4 |x_i|^2 \mu_i \mathcal{K}_0(2\|\mathbf{x}\boldsymbol{\mu}\|) & \text{if } \nu = 1 \\ \sigma^2 2\nu^\nu \frac{\Gamma(1-\nu)}{\Gamma(\nu)} \frac{1}{\mu_i} \left(\frac{|x_i| \mu_i}{\|\mathbf{x}\boldsymbol{\mu}\|^{1-\nu}} \right)^2 K_{1,1-\nu}(\|\mathbf{x}\boldsymbol{\mu}\|) & \text{if } \nu < 1. \end{cases} \quad (5.17)$$

Proof. The first assertion is a simple matter of differentiating Equation (5.7). In the following calculation, the fourth line is given by formula 9.6.28 (page 376) in Abramowitz and Stegun [1964].

$$\begin{aligned}
\frac{\partial}{\partial \mu_i} (\sigma^2 K_{r,\nu}(\mathbf{x}\boldsymbol{\mu})) &= \sigma^2 \frac{\partial}{\partial \mu_i} (K_{1,\nu}(\|\mathbf{x}\boldsymbol{\mu}\|)) \\
&= \sigma^2 x_i^2 \mu_i \|\mathbf{x}\boldsymbol{\mu}\|^{-1} K'_{1,\nu}(\|\mathbf{x}\boldsymbol{\mu}\|) \\
&= \sigma^2 x_i^2 \mu_i \|\mathbf{x}\boldsymbol{\mu}\|^{-1} \left(\frac{2\sqrt{\nu}}{\Gamma(\nu)2^{\nu-1}} \frac{d}{dy} \Big|_{y=2\sqrt{\nu}\|\mathbf{x}\boldsymbol{\mu}\|} [y^\nu \mathcal{K}_\nu(y)] \right) \\
&= \sigma^2 x_i^2 \mu_i \|\mathbf{x}\boldsymbol{\mu}\|^{-1} \left(\frac{2\sqrt{\nu}}{\Gamma(\nu)2^{\nu-1}} [-y \cdot y^{\nu-1} \mathcal{K}_{\nu-1}(y)]_{y=2\sqrt{\nu}\|\mathbf{x}\boldsymbol{\mu}\|} \right).
\end{aligned} \tag{5.18}$$

From there, Equation (5.16) follows immediately. Rewriting it in the form given in (5.17) only requires us to recall $\Gamma(\nu) = (\nu - 1)\Gamma(\nu - 1)$ (case $\nu > 1$), $\Gamma(1) = 1$ (case $\nu = 1$) and $\mathcal{K}_{\nu-1} = \mathcal{K}_{1-\nu}$ (case $\nu < 1$). \square

Accounting for low correlation: $\|\boldsymbol{\mu}\| \rightarrow \infty$

In this subsection, we consider a fixed design set of n coordinate-distinct points $\mathbf{x}^{(k)}$ ($k \in \llbracket 1, n \rrbracket$) in \mathbb{R}^r .

Lemma 5.7. *For any Matérn anisotropic geometric or tensorized correlation kernel with smoothness $\nu > 0$, for all $b < 2 \min(1, \nu) - 1$ and $c > 1$ (and if $\nu \neq 1$, for all $b \leq 2 \min(1, \nu) - 1$),*

- (a) $\forall \boldsymbol{\mu} \in (\mathbb{R}_+)^r$, $\|\frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}}\| \leq M_{i,1} \mu_i^{-c}$.
- (b) $\forall \boldsymbol{\mu} \in (\mathbb{R}_+)^r$, $\|\frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}}\| \leq M_{i,2} \mu_i^b$.

Proof. This can be gathered from Lemma 5.5 or 5.6 after recalling that 1) a Matérn kernel is a bounded function, 2) $\forall \nu \geq 0$, as $z \rightarrow +\infty$, $\mathcal{K}_\nu(z) \sim \sqrt{\pi} \exp(-z)/\sqrt{2z}$ (Abramowitz and Stegun [1964] 9.7.2) and 3) as $z \rightarrow 0$, $\mathcal{K}_0(z) \sim -\log(z)$ (Abramowitz and Stegun [1964] 9.6.8). \square

Let us define

$$f_i(\mu_i | \boldsymbol{\mu}_{-i}) := \sqrt{[\mathcal{I}(\boldsymbol{\mu})]_{ii}}; \tag{5.19}$$

$$\pi_i(\mu_i | \boldsymbol{\mu}_{-i}) := f_i(\mu_i | \boldsymbol{\mu}_{-i}) / \int_0^\infty f_i(\mu_i = t | \boldsymbol{\mu}_{-i}) dt. \tag{5.20}$$

Proposition 5.8. *For any Matérn anisotropic geometric or tensorized correlation kernel with smoothness $\nu > 0$, for all $\mu_i \in (0, +\infty)$, $\pi(\mu_i | \boldsymbol{\mu}_{-i})$, seen as a function of $\boldsymbol{\mu}$, is well defined and continuous over $\{\boldsymbol{\mu} \in [0, +\infty)^r : \mu_i \neq 0, \boldsymbol{\mu}_{-i} \neq \mathbf{0}_{r-1}\}$.*

Proof. For any given $\tilde{\boldsymbol{\mu}} \in [0, +\infty)^r$ such that $\tilde{\mu}_i \neq 0$ and $\tilde{\boldsymbol{\mu}}_{-i} \neq \mathbf{0}_{r-1}$, we prove that $\pi(\mu_i | \boldsymbol{\mu}_{-i})$, seen as a function of $\boldsymbol{\mu}$, is well defined and continuous at $\boldsymbol{\mu} = \tilde{\boldsymbol{\mu}}$.

For a start, notice that if $\boldsymbol{\mu}$ is confined to a sufficiently small neighborhood of $\tilde{\boldsymbol{\mu}}$, then $\|\boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{-1}\|$ remains bounded. Therefore, Lemma 5.7 implies that $\int_0^\infty f_i(\mu_i = t | \boldsymbol{\mu}_{-i}) dt$ is finite and, thanks to the dominated convergence theorem, that it is continuous at $\boldsymbol{\mu}_{-i} = \tilde{\boldsymbol{\mu}}_{-i}$. \square

Definition 5.9. An anisotropic geometric or tensorized correlation kernel is said to be “well-behaved” if its one-dimensional version is, for any set of parameters, a positive decreasing function on $[0, +\infty)$ that vanishes in the neighborhood of $+\infty$.

Lemma 5.10. Provided a coordinate-distinct design set is used, a well-behaved anisotropic geometric or tensorized correlation kernel parametrized by $\boldsymbol{\mu}$ has the following properties:

- (a) for any fixed $\boldsymbol{\mu}_{-i} \in [0, +\infty)^{r-1}$, it is a decreasing function of μ_i ;
 (b) as $\|\boldsymbol{\mu}\| \rightarrow \infty$, $\|\boldsymbol{\Sigma}_{\boldsymbol{\mu}} - \mathbf{I}_n\| \rightarrow 0$.

Lemma 5.11. For any well-behaved correlation kernel, as $\|\boldsymbol{\mu}\| \rightarrow \infty$, $\text{Tr} \left[\frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}} \boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{-1} \right] = o \left(\left\| \frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}} \right\| \right)$.

Proof. This result is due to the fact that all $\frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}}$'s diagonal coefficients are null and $\boldsymbol{\Sigma}_{\boldsymbol{\mu}}$ goes to the identity matrix as $\|\boldsymbol{\mu}\| \rightarrow \infty$. \square

Let us now define

$$h_i(\mu_i | \boldsymbol{\mu}_{-i}) := \sqrt{\text{Tr} \left[\left(\frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}} \right)^2 \right]} = \left\| \frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}} \right\|. \quad (5.21)$$

Lemma 5.12. For any well-behaved correlation kernel, as $\|\boldsymbol{\mu}\| \rightarrow \infty$, $f_i(\mu_i | \boldsymbol{\mu}_{-i}) \sim h_i(\mu_i | \boldsymbol{\mu}_{-i})$.

Proof. Because $\boldsymbol{\Sigma}_{\boldsymbol{\mu}}$ goes to the identity matrix, this is a direct consequence of Lemma 5.11. \square

Corollary 5.13. For any well-behaved correlation kernel, there exist $S > 0$, and $0 < a < b$ such that, whenever $\|\boldsymbol{\mu}\| \geq S$,

$$a h_i(\mu_i | \boldsymbol{\mu}_{-i}) \leq f_i(\mu_i | \boldsymbol{\mu}_{-i}) \leq b h_i(\mu_i | \boldsymbol{\mu}_{-i}). \quad (5.22)$$

In the following, $\boldsymbol{\Sigma}_{\boldsymbol{\mu}_{-i}}$ is the correlation matrix that would be obtained if μ_i were replaced by 0. Moreover, if \mathbf{M} is a matrix, $\mathbf{M}^{(kl)}$ is its element in the k -th row and l -th column.

Lemma 5.14. If a well-behaved correlation kernel is used, there exist real constants $S > 0$ and $c > 0$ such that, for all $\mu_i \in (0, +\infty)$ and whenever $\|\boldsymbol{\mu}_{-i}\| \geq S$,

$$\pi_i(\mu_i | \boldsymbol{\mu}_{-i}) \geq c \frac{\left\| \frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}} \right\|}{\sum_{k \neq l} \boldsymbol{\Sigma}_{\boldsymbol{\mu}_{-i}}^{(kl)}}. \quad (5.23)$$

Proof. If a well-behaved correlation kernel is used, then for any for any $\epsilon > 0$, Corollary 5.13 implies that

$$\int_0^{+\infty} f_i(\mu_i = t | \boldsymbol{\mu}_{-i}) dt \leq b \int_0^{+\infty} h_i(\mu_i = t | \boldsymbol{\mu}_{-i}) dt \leq -b \sum_{k \neq l} \int_0^{+\infty} \frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{(kl)} dt. \quad (5.24)$$

The last inequality holds because the Frobenius norm of any matrix is smaller than or equal to the sum of the absolute values of its elements and the correlation kernel is a decreasing function of μ_i . Now, for all $k \neq l$, when $\mu_i \rightarrow +\infty$, $\boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{(kl)} \rightarrow 0$ and when $\mu_i = 0$, $\boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{(kl)} = \boldsymbol{\Sigma}_{\boldsymbol{\mu}_{-i}}^{(kl)}$. From this, we gather that

$$\int_0^{+\infty} f_i(\mu_i = t | \boldsymbol{\mu}_{-i}) dt \leq b \sum_{k \neq l} \boldsymbol{\Sigma}_{\boldsymbol{\mu}_{-i}}^{(kl)}. \quad (5.25)$$

From this, we deduce that

$$\pi_i(\mu_i | \boldsymbol{\mu}_{-i}) = \frac{f_i(\mu_i | \boldsymbol{\mu}_{-i})}{\int_0^{+\infty} f_i(\mu_i = t | \boldsymbol{\mu}_{-i}) dt} \geq \frac{a}{b} \frac{\|\frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}}\|}{\sum_{k \neq l} \boldsymbol{\Sigma}_{\boldsymbol{\mu}_{-i}}^{(kl)}}. \quad (5.26)$$

□

This Lemma has the following immediate consequence:

Proposition 5.15. *If a well-behaved tensorized kernel is used, there exists $S > 0$ and for every $i \in \llbracket 1, r \rrbracket$, there exists a function $M_i : (0, +\infty) \rightarrow (0, +\infty)$ such that for all $\|\boldsymbol{\mu}_{-i}\| \geq S$, $\pi_i(\mu_i | \boldsymbol{\mu}_{-i}) \geq M_i(\mu_i)$.*

Proof. If a tensorized correlation kernel is used, for every pair of integers $(k, l) \in \llbracket 1, r \rrbracket^2$ such that $k \neq l$, define the function $M_i^{(kl)} : (0, +\infty) \rightarrow (0, +\infty)$; $t \mapsto \left| \frac{d}{dt} \boldsymbol{\Sigma}_{\mu_i=t, \boldsymbol{\mu}_{-i}=\mathbf{0}_{r-1}}^{(kl)} \right|$.

$$\left\| \frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}} \right\| \geq \frac{1}{n} \sum_{k \neq l} \left| \frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{(kl)} \right| = \frac{1}{n} \sum_{k \neq l} M_i^{(kl)}(\mu_i) \boldsymbol{\Sigma}_{\boldsymbol{\mu}_{-i}}^{(kl)} \geq \frac{1}{n} \min_{k \neq l} M_i^{(kl)}(\mu_i) \sum_{k \neq l} \boldsymbol{\Sigma}_{\boldsymbol{\mu}_{-i}}^{(kl)}. \quad (5.27)$$

This fact, joined with Lemma 5.14, yields the result. □

Proposition 5.16. *Assume a well-behaved anisotropic geometric correlation kernel is used. If the corresponding one-dimensional kernel K has the properties (P1) and (P2), then for every $i \in \llbracket 1, r \rrbracket$, there exist positive functions s_i and m_i defined on $(0, +\infty)$ such that, for all $\|\boldsymbol{\mu}_{-i}\| \geq s_i(\mu_i)$, $\pi_i(\mu_i | \boldsymbol{\mu}_{-i}) \geq m_i(\mu_i)$.*

(P1) : *There exist $S_1 > 0$ and $M_1 > 0$ such that, for all $t \geq S_1$, $|K'(t)| \geq M_1 t K(t)$.*

(P2) : *For any $a > 0$, there exist $S_2(a) > 0$ and $M_2(a) > 0$ such that, whenever $t \geq S_2(a)$, $K(t+a) \geq M_2(a)K(t)$.*

Proof. From (P1), we gather that for all $a > 0$ and for all $t \geq S_1$, $|K'(\sqrt{t^2 + a^2})| \geq M_1 \sqrt{t^2 + a^2} K(\sqrt{t^2 + a^2})$. Now, because the correlation kernel is well-behaved, K is a decreasing function. As $\sqrt{t^2 + a^2} \leq t + a$, $K(\sqrt{t^2 + a^2}) \geq K(t + a)$.

Plugging this into the previous inequality, we get $|K'(\sqrt{t^2 + a^2})| \geq M_1 \sqrt{t^2 + a^2} K(t + a)$.

If $t \geq \max(S_1, S_2(a))$, we can then use (P2) to obtain

$$|K'(\sqrt{t^2 + a^2})| \geq M_1 M_2(a) \sqrt{t^2 + a^2} K(t). \quad (5.28)$$

Independently from this, we have the following algebraic fact:

$$\left\| \frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}} \right\| \geq \frac{1}{n} \sum_{k \neq l} \left| \frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{(kl)} \right|. \quad (5.29)$$

Because we use a well-behaved anisotropic geometric kernel, defining the function $M_i^{(kl)} : (0, +\infty) \rightarrow (0, +\infty)$; $t \mapsto \left(x_j^{(k)} - x_j^{(l)} \right)^2$, we can write:

$$\left| \frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{(kl)} \right| = - \frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{(kl)} = \left(x_i^{(k)} - x_i^{(l)} \right)^2 \mu_i \frac{K'(\|(\mathbf{x}^{(k)} - \mathbf{x}^{(l)})\boldsymbol{\mu}\|)}{\|(\mathbf{x}^{(k)} - \mathbf{x}^{(l)})\boldsymbol{\mu}\|}. \quad (5.30)$$

Setting $a_{kl} := |x_i^{(k)} - x_i^{(l)}| \mu_i$ and $t_{kl} := \|(\mathbf{x}_{-i}^{(k)} - \mathbf{x}_{-i}^{(l)})\boldsymbol{\mu}_{-i}\|$ (and thus, naturally, $\sqrt{t_{kl}^2 + a_{kl}^2} = \|(\mathbf{x}^{(k)} - \mathbf{x}^{(l)})\boldsymbol{\mu}\|$), and provided $\|\boldsymbol{\mu}_{-i}\|$ is sufficiently large to make all t_{kl} s meet the conditions

necessary to apply (P1) and (P2) (that depend in the case of (P2) on the a_{kl} s), Equation (5.28) yields the existence of some number $m_i^{(kl)}(\mu_i) > 0$ such that

$$\frac{K'(\|(\mathbf{x}^{(k)} - \mathbf{x}^{(l)})\boldsymbol{\mu}\|)}{\|(\mathbf{x}^{(k)} - \mathbf{x}^{(l)})\boldsymbol{\mu}\|} \geq m_i^{(kl)}(\mu_i)K(\|(\mathbf{x}_{-i}^{(k)} - \mathbf{x}_{-i}^{(l)})\boldsymbol{\mu}_{-i}\|) = m_i^{(kl)}(\mu_i)\boldsymbol{\Sigma}_{\boldsymbol{\mu}_{-i}}^{(kl)}. \quad (5.31)$$

Finally, setting $m_i(\mu_i) := \mu_i \min_{k \neq l} \left[\left(x_i^{(k)} - x_i^{(l)} \right)^2 m_i^{(kl)}(\mu_i) \right]$, we get

$$\left\| \frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}} \right\| \geq \frac{m_i(\mu_i)}{n} \sum_{k \neq l} \boldsymbol{\Sigma}_{\boldsymbol{\mu}_{-i}}^{(kl)}. \quad (5.32)$$

Then, applying Lemma 5.14 yields the result. \square

Proposition 5.17. *Matérn one-dimensional kernels with smoothness parameter $\nu > 1$ have the properties (P1) and (P2) of Proposition 5.16.*

Proof. (P1) is given by Lemma 5.6, after noticing that, denoting by K_ν the Matérn one-dimensional kernel of smoothness $\nu > 1$, provided t is sufficiently large, $K_\nu(t) \leq K_{\nu-1}(t)$.

This inequality ensues from the fact that $\forall \nu \geq 0$, as $t \rightarrow +\infty$, $\mathcal{K}_\nu(t) \sim \sqrt{\pi} \exp(-t)/\sqrt{2t}$ (Abramowitz and Stegun [1964] 9.7.2), hence $K_\nu(t) \sim 2/\Gamma(\nu)(\sqrt{\nu t})^\nu \sqrt{\pi/(4\sqrt{\nu t})} \exp(-2\sqrt{\nu t})$. Moreover, this last equivalence relation also implies (P2). \square

Proposition 5.18. *For Matérn anisotropic geometric kernels with smoothness $\nu > 1$ and Matérn tensorized correlation kernels with smoothness $\nu > 0$, for any $\delta > 0$, $i \in \llbracket 1, r \rrbracket$ and $\mu_i \in (0, +\infty)$, there exists $b_{i,\delta}(\mu_i) > 0$ such that, if $\|\boldsymbol{\mu}_{-i}\| \geq \delta$, then $\pi_i(\mu_i|\boldsymbol{\mu}_{-i}) \geq b_{i,\delta}(\mu_i)$.*

Proof. Matérn correlation kernels with such smoothness parameters make Proposition 5.15 or 5.16 applicable. Therefore, there exist $s_i(\mu_i) > 0$ and $m_i(\mu_i) > 0$ such that, if $\|\boldsymbol{\mu}_{-i}\| \geq s_i(\mu_i)$, $\pi_i(\mu_i|\boldsymbol{\mu}_{-i}) \geq m_i(\mu_i)$. Besides, we know from Proposition 5.8 that $\pi_i(\mu_i|\boldsymbol{\mu}_{-i})$, seen as a function of $\boldsymbol{\mu}_{-i}$, is continuous and positive over the compact set $\{\boldsymbol{\mu}_{-i} : \delta \leq \|\boldsymbol{\mu}_{-i}\| \leq s_i(\mu_i)\}$. Thus its minimum $\tilde{m}_{i,\delta}(\mu_i)$ on this set is positive and we obtain the result by setting $b_{i,\delta}(\mu_i) := \min(m_i(\mu_i), \tilde{m}_{i,\delta}(\mu_i))$. \square

Proposition 5.19. *For Matérn anisotropic geometric correlation kernels with smoothness $\nu > 1$ and for Matérn tensorized correlation kernels with smoothness $\nu > 0$, for any $\mathbf{y} \in \mathbb{R}^n \setminus \{0\}^n$, any $\delta > 0$ and any $\mu_i \in (0, +\infty)$, there exists $b_{i,\delta,\mathbf{y}}(\mu_i) > 0$ such that, if $\|\boldsymbol{\mu}_{-i}\| \geq \delta$, then $\pi_i(\mu_i|\mathbf{y}, \boldsymbol{\mu}_{-i}) \geq b_{i,\delta,\mathbf{y}}(\mu_i)$.*

Proof. Set $\delta > 0$ and $\mathbf{y} \in \mathbb{R}^n \setminus \{0\}^n$. There exist $m_\delta > 0$ and $M_\delta > 0$ s.t. $\forall \boldsymbol{\mu} \in (0, +\infty)^r$, $\|\boldsymbol{\mu}_{-i}\| \geq \delta \Rightarrow m_\delta \leq \|\boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{-1}\| \leq M_\delta$, so there also exist $m_{\delta,\mathbf{y}} > 0$ and $M_{\delta,\mathbf{y}} > 0$ s.t. $m_{\delta,\mathbf{y}} \leq L(\mathbf{y}|\boldsymbol{\mu}) \leq M_{\delta,\mathbf{y}}$. This, combined with Proposition 5.18, yields the result. \square

Accounting for high correlation: $\|\boldsymbol{\mu}\| \rightarrow 0$

This part of the proof relies on the combination of some spectral study of the Matérn kernels and on the study of the matrices that are part of the series expansion of the correlation matrix $\boldsymbol{\Sigma}_{\boldsymbol{\mu}}$ when $\|\boldsymbol{\mu}\| \rightarrow 0$ for three types of Matérn kernels: isotropic, tensorized and anisotropic geometric.

Lemma 5.20. *There exists a covariance kernel $\tilde{K}_{r,\nu}$ such that for any design set $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$, for all $\boldsymbol{\mu} \in (\mathbb{R}_+)^r$ and $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n) \in \mathbb{R}^n$,*

$$\begin{aligned} & \sum_{j,k=1}^n \xi_j \xi_k K_{r,\nu} \left(\left(\mathbf{x}^{(j)} - \mathbf{x}^{(k)} \right) \boldsymbol{\mu} \right) \\ & \geq 2^{-\frac{r}{2}-\nu} M_r(\nu) f_{r,\nu}(\|\boldsymbol{\mu}\|_\infty) \sum_{j,k=1}^n \xi_j \xi_k \tilde{K}_{r,\nu} \left(\frac{\boldsymbol{\mu}}{\|\boldsymbol{\mu}\|_\infty} \left(\mathbf{x}^{(j)} - \mathbf{x}^{(k)} \right) \right), \end{aligned} \quad (5.33)$$

where $f_{r,\nu}(t) = (2\sqrt{\nu})^{-r-2\nu} t^{-r}$ if $t \geq (2\sqrt{\nu})^{-1}$ and $f_{r,\nu}(t) = t^{2\nu}$ if $t \leq (2\sqrt{\nu})^{-1}$.

Proof. For all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^r$, $K_{r,\nu}(\mathbf{x} - \mathbf{y}) = \int_{\mathbb{R}^r} \widehat{K}_{r,\nu}(\boldsymbol{\omega}) e^{i\langle \boldsymbol{\omega}, \mathbf{x} - \mathbf{y} \rangle} d\boldsymbol{\omega}$.

$$\begin{aligned} & \sum_{j,k=1}^n \xi_j \xi_k K_{r,\nu} \left(\left(\mathbf{x}^{(j)} - \mathbf{x}^{(k)} \right) \boldsymbol{\mu} \right) \\ & = \int_{\mathbb{R}^r} \widehat{K}_{r,\nu}(\boldsymbol{\omega}) \left| \sum_{j=1}^n \xi_j e^{i\langle \boldsymbol{\omega}, \mathbf{x}^{(j)} \rangle} \right|^2 d\boldsymbol{\omega} \\ & = M_r(\nu) \|\boldsymbol{\mu}\|_\infty^{-r} \int_{\mathbb{R}^r} (4\nu + \|\boldsymbol{\mu}\|_\infty^{-2} \|\mathbf{s}\|^2)^{-\frac{r}{2}-\nu} \left| \sum_{j=1}^n \xi_j e^{i\langle \frac{\boldsymbol{\mu}}{\|\boldsymbol{\mu}\|_\infty} \mathbf{s}, \mathbf{x}^{(j)} \rangle} \right|^2 d\mathbf{s} \\ & \geq 2^{-\frac{r}{2}-\nu} M_r(\nu) f_{r,\nu}(\|\boldsymbol{\mu}\|_\infty) \int_{\mathbb{R}^r \setminus B(0,1)} \|\mathbf{s}\|^{-r-2\nu} \left| \sum_{j=1}^n \xi_j e^{i\langle \mathbf{s}, \frac{\boldsymbol{\mu}}{\|\boldsymbol{\mu}\|_\infty} \mathbf{x}^{(j)} \rangle} \right|^2 d\mathbf{s}. \end{aligned} \quad (5.34)$$

Now, let $\tilde{K}_{r,\nu}$ be the function with Fourier transform $\widehat{K}_{r,\nu}(\boldsymbol{\omega}) = \mathbf{1}_{\{\|\boldsymbol{\omega}\| \geq 1\}} \|\boldsymbol{\omega}\|^{-r-2\nu}$. According to Bochner's theorem, $\tilde{K}_{r,\nu}$ is a correlation kernel, which leads to the conclusion. \square

Lemma 5.21. *For every design set with coordinate-distinct points $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$, there exists a constant $c_x > 0$ such that for all $\boldsymbol{\mu} \in (\mathbb{R}_+)^r$,*

$$\forall \boldsymbol{\xi} = (\xi_1, \dots, \xi_n) \in \mathbb{R}^n, \quad \sum_{j,k=1}^n \xi_j \xi_k K_{r,\nu} \left(\left(\mathbf{x}^{(j)} - \mathbf{x}^{(k)} \right) \boldsymbol{\mu} \right) \geq c_x \|\boldsymbol{\xi}\|^2 2^{-\frac{r}{2}-\nu} M_r(\nu) f_{r,\nu}(\|\boldsymbol{\mu}\|_\infty) \quad (5.35)$$

where $f_{r,\nu}(t) = (2\sqrt{\nu})^{-r-2\nu} t^{-r}$ if $t \geq (2\sqrt{\nu})^{-1}$ and $f_{r,\nu}(t) = t^{2\nu}$ if $t \leq (2\sqrt{\nu})^{-1}$.

Proof. For every design set $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$, the set of all design sets that can be written $\frac{\boldsymbol{\mu}}{\|\boldsymbol{\mu}\|_\infty} \mathbf{x}^{(1)}, \dots, \frac{\boldsymbol{\mu}}{\|\boldsymbol{\mu}\|_\infty} \mathbf{x}^{(n)}$ ($\boldsymbol{\mu} \in (\mathbb{R}_+)^r$) is compact. If the design set $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$ has coordinate-distinct points, then every design set in the aforementioned compact set has no overlapping points. Thus the conclusion follows from Lemma 5.20. \square

Proposition 5.22. *With Matérn anisotropic geometric or tensorized kernels, for every design set with coordinate-distinct points $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$, as $\|\boldsymbol{\mu}\| \rightarrow 0$, $\|\boldsymbol{\Sigma}_\boldsymbol{\mu}^{-1}\| = O(\|\boldsymbol{\mu}\|^{-2\nu})$.*

Proof. For Matérn anisotropic geometric kernels, we need only apply Lemma 5.21. In the case of tensorized Matérn kernels, analogous results to Lemma 5.20 and then Lemma 5.21 may be used. \square

Abramowitz and Stegun [1964] give the following results on the modified Bessel function of second kind (usually noted K_ν and which we note \mathcal{K}_ν in order to avoid confusion with the

Matérn correlation kernel). If I_ν is the modified Bessel function of first kind and ψ is the function defined in (6.3.2) by $\psi : \mathbb{N} \setminus \{0\} \rightarrow \mathbb{R} ; k \mapsto -\gamma + \sum_{i=1}^{k-1} i^{-1}$:

$$I_\nu(z) = \left(\frac{1}{2}z\right)^\nu \sum_{k=0}^{\infty} \frac{\left(\frac{1}{4}z^2\right)^k}{k!\Gamma(\nu+k+1)} \quad (9.6.10)$$

$$\mathcal{K}_\nu(z) = \frac{1}{2\pi} \frac{I_{-\nu}(z) - I_\nu(z)}{\sin(\nu z)} \quad \text{if } \nu \notin \mathbb{Z}. \quad (9.6.2)$$

This gives us the series expansion of $K_{1,\nu}(z)$ ($\nu \in [0, +\infty) \setminus \mathbb{N}$) when $z \rightarrow 0$:

$$\begin{aligned} & K_{1,\nu}(z) \\ &= \frac{\pi}{\Gamma(\nu) \sin(\nu\pi)} \left(\sum_{0 \leq k < \nu} \frac{\nu^k z^{2k}}{k!\Gamma(-\nu+k+1)} - \frac{\nu^\nu z^{2\nu}}{\Gamma(\nu+1)} + o(z^{2\nu}) \right) \\ &= \frac{\pi}{\Gamma(\nu) \sin(\nu\pi)\Gamma(-\nu+1)} \left(\sum_{0 \leq k < \nu} \frac{\Gamma(-\nu+1)}{k!\Gamma(-\nu+k+1)} \nu^k z^{2k} - \frac{\Gamma(-\nu+1)}{\Gamma(\nu+1)} \nu^\nu z^{2\nu} + o(z^{2\nu}) \right) \\ &= \sum_{0 \leq k < \nu} \frac{\Gamma(-\nu+1)}{k!\Gamma(-\nu+k+1)} \nu^k z^{2k} + \frac{\Gamma(-\nu)}{\Gamma(\nu)} \nu^\nu z^{2\nu} + o(z^{2\nu}) \\ &= \sum_{0 \leq k < \nu} (-1)^k \frac{\Gamma(\nu-k)}{k!\Gamma(\nu)} \nu^k z^{2k} + \frac{\Gamma(-\nu)}{\Gamma(\nu)} \nu^\nu z^{2\nu} + o(z^{2\nu}). \end{aligned} \quad (5.36)$$

In the remainder of this subsection, we consider a fixed design set with n coordinate-distinct points $\mathbf{x}^{(k)}$ ($k \in \llbracket 1, n \rrbracket$) in \mathbb{R}^r . Moreover, all Matérn kernels we consider are assumed to have non-integer smoothness parameter ν .

Let us now define, for every nonnegative integer $k < \nu$ the matrix \mathbf{D}^k whose (i, j) element is

$$\mathbf{D}^k(i, j) := (-1)^k \frac{\Gamma(\nu-k)}{k!\Gamma(\nu)} \nu^k \left\| \mathbf{x}^{(j)} - \mathbf{x}^{(k)} \right\|^{2k}. \quad (5.37)$$

Let us also define the matrix \mathbf{D}^ν whose (i, j) element is

$$\mathbf{D}^\nu(i, j) := \frac{\Gamma(-\nu)}{\Gamma(\nu)} \nu^\nu \left\| \mathbf{x}^{(j)} - \mathbf{x}^{(k)} \right\|^{2\nu} \quad \text{if } \nu \in [0, +\infty) \setminus \mathbb{N}. \quad (5.38)$$

If the correlation kernel is Matérn isotropic, Σ_μ has the following series expansion if ν is not an integer when $\mu \rightarrow 0+$:

$$\Sigma_\mu = \sum_{0 \leq k < \nu} \mu^{2k} \mathbf{D}^k + \mu^{2\nu} \mathbf{D}^\nu + \mathbf{R}_\mu. \quad (5.39)$$

In this expansion, $\mu^{-2\nu} \|\mathbf{R}_\mu\| \rightarrow 0$.

For any integer $i \in \llbracket 1, r \rrbracket$ and any nonnegative integer $k < \nu$ define the matrix \mathbf{D}_i^k whose (m, p) element is

$$\mathbf{D}_i^k(m, p) := (-1)^k \frac{\Gamma(\nu-k)}{k!\Gamma(\nu)} \nu^k \left| x_i^{(m)} - x_i^{(p)} \right|^{2k} \quad (5.40)$$

and also the matrix \mathbf{D}_i^ν whose (m, p) element is

$$\mathbf{D}_i^\nu(m, p) := \frac{\Gamma(-\nu)}{\Gamma(\nu)} \nu^\nu \left| x_i^{(m)} - x_i^{(p)} \right|^{2\nu}. \quad (5.41)$$

For every $i \in \llbracket 1, r \rrbracket$, if the points in the design set differed only through their i -th coordinate, the series expansion of the correlation matrix (using a Matérn anisotropic geometric or tensorized kernel) when $\|\boldsymbol{\mu}\| \rightarrow 0$ (and thus when $\mu_i \rightarrow 0$) would be

$$\boldsymbol{\Sigma}_{\mu_i} = \sum_{0 \leq k < \nu} \mu_i^{2k} \mathbf{D}_i^k + \mu_i^{2\nu} \mathbf{D}_i^\nu + \mathbf{R}_{\mu_i} \quad (5.42)$$

where $\mu_i^{-2\nu} \|\mathbf{R}_{\mu_i}\| \rightarrow 0$ as $\mu_i \rightarrow 0$.

Note the following identities:

$$\mathbf{D}_i^0 = \mathbf{1}\mathbf{1}^\top; \quad (5.43)$$

$$\mathbf{D}_i^1 = -\frac{\Gamma(\nu-1)}{\Gamma(\nu)} \nu \left\{ \mathbf{1} (\mathbf{X}_i^{\circ 2})^\top + (\mathbf{X}_i^{\circ 2}) \mathbf{1}^\top - 2\mathbf{X}_i \mathbf{X}_i^\top \right\}; \quad (5.44)$$

$$\begin{aligned} \mathbf{D}_i^2 = \frac{\Gamma(\nu-2)}{\Gamma(\nu)} \nu^2 \left\{ \mathbf{1} (\mathbf{X}_i^{\circ 4})^\top + (\mathbf{X}_i^{\circ 4}) \mathbf{1}^\top - 4\mathbf{X}_i (\mathbf{X}_i^{\circ 3})^\top \right. \\ \left. - 4(\mathbf{X}_i^{\circ 3}) \mathbf{X}_i^\top + 6(\mathbf{X}_i^{\circ 2}) (\mathbf{X}_i^{\circ 2})^\top \right\}. \end{aligned} \quad (5.45)$$

If a tensorized correlation kernel is used, the correlation matrix $\boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{tens}$ may be written

$$\boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{tens} = \overset{\circ}{\prod}_{i \in \llbracket 1, r \rrbracket} \boldsymbol{\Sigma}_{\mu_i} \quad (5.46)$$

where the subscript \circ above the symbol \prod serves to denote the Hadamard product of matrices.

In case a Matérn anisotropic geometric kernel is used, then define for any nonnegative interger $k < \nu$ the matrix $\mathbf{D}^k(\boldsymbol{\mu})$ whose (m, p) element is

$$\mathbf{D}^k(\boldsymbol{\mu})(m, p) := (-1)^k \frac{\Gamma(\nu-k)}{k! \Gamma(\nu)} \nu^k d_{m,p}(\boldsymbol{\mu})^{2k} \quad (5.47)$$

where $d_{m,p}(\boldsymbol{\mu}) = \|(\mathbf{x}^{(m)} - \mathbf{x}^{(p)}) \boldsymbol{\mu}\|$.

And, similarly, we may define the matrix $\mathbf{D}^\nu(\boldsymbol{\mu})$ whose (m, p) element is

$$\mathbf{D}^\nu(\boldsymbol{\mu})(m, p) := \frac{\Gamma(-\nu)}{\Gamma(\nu)} \nu^\nu d_{m,p}(\boldsymbol{\mu})^{2\nu} \quad \text{if } \nu \in [0, +\infty) \setminus \mathbb{N}. \quad (5.48)$$

We thus have (if $\nu \in [0, +\infty) \setminus \mathbb{N}$)

$$\boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{geom} = \sum_{0 \leq k < \nu} \mathbf{D}^k(\boldsymbol{\mu}) + \mathbf{D}^\nu(\boldsymbol{\mu}) + \mathbf{R}_{\boldsymbol{\mu}}^{geom} \quad (5.49)$$

where $\|\boldsymbol{\mu}\|^{-2\nu} \|\mathbf{R}_{\boldsymbol{\mu}}^{geom}\| \rightarrow 0$ as $\|\boldsymbol{\mu}\| \rightarrow 0$.

Similar identities to those of Equation (5.43) can be derived to make Equation (5.49) more explicit for small values of ν .

$$\mathbf{D}^0(\boldsymbol{\mu}) = \mathbf{1}\mathbf{1}^\top; \quad (5.50)$$

$$\mathbf{D}^1(\boldsymbol{\mu}) = -\frac{\nu}{\nu-1} \left\{ \sum_{i=1}^r \mu_i^2 \left(\mathbf{1} (\mathbf{X}_i^{\circ 2})^\top + (\mathbf{X}_i^{\circ 2}) \mathbf{1}^\top - 2\mathbf{X}_i \mathbf{X}_i^\top \right) \right\}; \quad (5.51)$$

$$\begin{aligned} \mathbf{D}^2(\boldsymbol{\mu}) = \frac{\nu^2}{(\nu-1)(\nu-2)} \left\{ \sum_{i,j \in \llbracket 1, r \rrbracket} \mu_i^2 \mu_j^2 \left(\mathbf{1} (\mathbf{X}_i^{\circ 2} \circ \mathbf{X}_j^{\circ 2})^\top + (\mathbf{X}_i^{\circ 2} \circ \mathbf{X}_j^{\circ 2}) \mathbf{1}^\top \right. \right. \\ \left. \left. - 2\mathbf{X}_i (\mathbf{X}_i \circ \mathbf{X}_j^{\circ 2})^\top - 2(\mathbf{X}_i \circ \mathbf{X}_j^{\circ 2}) \mathbf{X}_i^\top - 2\mathbf{X}_j (\mathbf{X}_j \circ \mathbf{X}_i^{\circ 2})^\top - 2(\mathbf{X}_j \circ \mathbf{X}_i^{\circ 2}) \mathbf{X}_j^\top \right. \right. \\ \left. \left. + (\mathbf{X}_i^{\circ 2}) (\mathbf{X}_j^{\circ 2})^\top + (\mathbf{X}_j^{\circ 2}) (\mathbf{X}_i^{\circ 2})^\top + 4(\mathbf{X}_i \circ \mathbf{X}_j) (\mathbf{X}_i \circ \mathbf{X}_j)^\top \right) \right\}. \quad (5.52) \end{aligned}$$

Fortunately, for small values of ν , $\boldsymbol{\Sigma}_\mu^{tens}$ can also be simply written.

$$\text{For } \nu \in (0, 1): \quad \boldsymbol{\Sigma}_\mu^{tens} = \mathbf{1}\mathbf{1}^\top + \sum_{i=1}^r \mu_i^{2\nu} \mathbf{D}_i^\nu + \mathbf{R}_\mu^{tens}. \quad (5.53)$$

$$\text{For } \nu \in (1, 2): \quad \boldsymbol{\Sigma}_\mu^{tens} = \mathbf{1}\mathbf{1}^\top + \mathbf{D}^1(\boldsymbol{\mu}) + \sum_{i=1}^r \mu_i^{2\nu} \mathbf{D}_i^\nu + \mathbf{R}_\mu^{tens}. \quad (5.54)$$

$$\text{For } \nu \in (2, 3): \quad \boldsymbol{\Sigma}_\mu^{tens} = \mathbf{1}\mathbf{1}^\top + \mathbf{D}^1(\boldsymbol{\mu}) + \frac{\nu-2}{\nu-1} \mathbf{D}^2(\boldsymbol{\mu}) + \sum_{i=1}^r \mu_i^4 \mathbf{D}_i^2 + \sum_{i=1}^r \mu_i^{2\nu} \mathbf{D}_i^\nu + \mathbf{R}_\mu^{tens}. \quad (5.55)$$

In the three expressions above, $\|\boldsymbol{\mu}\|^{-2\nu} \|\mathbf{R}_\mu^{tens}\| \rightarrow 0$ as $\|\boldsymbol{\mu}\| \rightarrow 0$.

Define k_ν as the orthogonal complement in \mathbb{R}^n of the vector space spanned by:

1. if $\nu \in (0, 1)$: $\mathbf{1}$;
2. if $\nu \in (1, 2)$: $\mathbf{1}$ and \mathbf{X}_i ($i \in \llbracket 1, r \rrbracket$);
3. if $\nu \in (2, 3)$: $\mathbf{1}$ and \mathbf{X}_i ($i \in \llbracket 1, r \rrbracket$) and $\mathbf{X}_i \circ \mathbf{X}_j$ ($i, j \in \llbracket 1, r \rrbracket$).

Clearly, for any $\nu \in (0, 1) \cup (1, 2) \cup (2, 3)$, for any vector $\mathbf{v} \in k_\nu$,

$$\mathbf{v}^\top \boldsymbol{\Sigma}_\mu^{geom} \mathbf{v} = \mathbf{v}^\top \mathbf{D}^\nu(\boldsymbol{\mu}) \mathbf{v} + \mathbf{v}^\top \mathbf{R}_\mu^{geom} \mathbf{v}, \quad (5.56)$$

$$\mathbf{v}^\top \boldsymbol{\Sigma}_\mu^{tens} \mathbf{v} = \sum_{i=1}^r \mu_i^{2\nu} \mathbf{v} \mathbf{D}_i^\nu \mathbf{v} + \mathbf{v}^\top \mathbf{R}_\mu^{tens} \mathbf{v}. \quad (5.57)$$

Since when $\mu \rightarrow 0$ $\|\mathbf{D}^\nu(\boldsymbol{\mu})\| = O(\|\boldsymbol{\mu}\|^{2\nu})$, $\|\mathbf{R}_\mu^{geom}\| = o(\|\boldsymbol{\mu}\|^{2\nu})$ and $\|\mathbf{R}_\mu^{tens}\| = o(\|\boldsymbol{\mu}\|^{2\nu})$, for any $\boldsymbol{\mu} \in (0, +\infty)^r$ such that $\|\boldsymbol{\mu}\|$ is small enough, there exists $c > 0$ such that for any $\mathbf{v} \in k_\nu$,

$$\max(\mathbf{v}^\top \boldsymbol{\Sigma}_\mu^{geom} \mathbf{v}, \mathbf{v}^\top \boldsymbol{\Sigma}_\mu^{tens} \mathbf{v}) \leq c \|\boldsymbol{\mu}\|^{2\nu} \mathbf{v}^\top \mathbf{v}. \quad (5.58)$$

Proposition 5.23. *For a Matérn anisotropic geometric or tensorized correlation kernel with smoothness parameter $\nu \in (0, 1) \cup (1, 2) \cup (2, 3)$, for any vector $\mathbf{y} \in \mathbb{R}^n$ not orthogonal to k_ν , when $\|\boldsymbol{\mu}\| \rightarrow 0$, $\|\boldsymbol{\mu}\|^{-2\nu} = O(\mathbf{y}^\top \boldsymbol{\Sigma}_\mu^{-1} \mathbf{y})$.*

Proof. Let d_ν be the dimension of k_ν and let \mathbf{O}_ν be an orthogonal $n \times n$ matrix whose first $n - d_\nu$ columns form an orthonormal basis of k_ν^\perp and whose last d_ν columns form an orthonormal basis of k_ν .

Then $\Sigma_\mu = \mathbf{O}_\nu \mathbf{O}_\nu^\top \Sigma_\mu \mathbf{O}_\nu \mathbf{O}_\nu^\top$. Consider the following decomposition of $\mathbf{O}_\nu^\top \Sigma_\mu \mathbf{O}_\nu$:

$$\mathbf{O}_\nu^\top \Sigma_\mu \mathbf{O}_\nu = \begin{pmatrix} \mathbf{A}_\mu & \mathbf{B}_\mu \\ \mathbf{B}_\mu^\top & \mathbf{C}_\mu \end{pmatrix} \quad (5.59)$$

where the blocks \mathbf{A}_μ , \mathbf{B}_μ and \mathbf{C}_μ are respectively $(n - d_\nu) \times (n - d_\nu)$, $(n - d_\nu) \times d_\mu$ and $d_\mu \times d_\mu$ matrices. Note that \mathbf{A}_μ and \mathbf{C}_μ represent the restriction of the scalar product defined by Σ_μ to k_ν^\perp and k_ν respectively. When $\|\boldsymbol{\mu}\|$ is small enough, defining $c > 0$ as in Equation (5.58), $\|\mathbf{C}_\mu\| \leq c\|\boldsymbol{\mu}\|^{2\nu}$.

$$\mathbf{O}_\nu^\top \Sigma_\mu^{-1} \mathbf{O}_\nu = \begin{pmatrix} \mathbf{I}_{n-d_\nu} & \mathbf{0} \\ -\mathbf{B}_\mu \mathbf{C}_\mu^{-1} & \mathbf{I}_{d_\nu} \end{pmatrix} \begin{pmatrix} (\mathbf{A}_\mu - \mathbf{B}_\mu \mathbf{C}_\mu^{-1} \mathbf{B}_\mu^\top)^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_\mu^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{I}_{n-d_\nu} & -\mathbf{C}_\mu^{-1} \mathbf{B}_\mu^\top \\ \mathbf{0} & \mathbf{I}_{d_\nu} \end{pmatrix}. \quad (5.60)$$

For any vector $\mathbf{y} \in \mathbb{R}^n$, there exist $\mathbf{y}_1 \in \mathbb{R}^{n-d_\nu}$ and $\mathbf{y}_2 \in \mathbb{R}^{d_\nu}$ such that

$$\mathbf{O}_\nu^\top \mathbf{y} = \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{pmatrix}, \quad (5.61)$$

$$\mathbf{y}^\top \Sigma_\mu^{-1} \mathbf{y} = \begin{pmatrix} \mathbf{y}_1 - \mathbf{C}_\mu^{-1} \mathbf{B}_\mu \mathbf{y}_2 \\ \mathbf{y}_2 \end{pmatrix}^\top \begin{pmatrix} (\mathbf{A}_\mu - \mathbf{B}_\mu \mathbf{C}_\mu^{-1} \mathbf{B}_\mu^\top)^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_\mu^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{y}_1 - \mathbf{C}_\mu^{-1} \mathbf{B}_\mu \mathbf{y}_2 \\ \mathbf{y}_2 \end{pmatrix}. \quad (5.62)$$

Given Σ_μ^{-1} is positive definite, the diagonal block $(\mathbf{A}_\mu - \mathbf{B}_\mu \mathbf{C}_\mu^{-1} \mathbf{B}_\mu^\top)^{-1}$ is positive definite too. This implies $\mathbf{y}^\top \Sigma_\mu^{-1} \mathbf{y} \geq \mathbf{y}_2^\top \mathbf{C}_\mu^{-1} \mathbf{y}_2$. When $\|\boldsymbol{\mu}\|$ is small enough, $\mathbf{y}_2^\top \mathbf{C}_\mu^{-1} \mathbf{y}_2 \geq c^{-1} \|\boldsymbol{\mu}\|^{-2\nu} \|\mathbf{y}_2\|^2$.

If \mathbf{y} is not orthogonal to k_ν , then $\|\mathbf{y}_2\| \neq 0$ and thus $\|\boldsymbol{\mu}\|^{-2\nu} = O(\mathbf{y}^\top \Sigma_\mu^{-1} \mathbf{y})$. □

Proposition 5.24. *Assume $\nu \in (1, 2) \cup (2, 3)$. For every $\mathbf{y} \in \mathbb{R}^n$ that is not orthogonal to the vector subspace k_ν , $L(\mathbf{y}|\boldsymbol{\mu}) f_i(\mu_i|\boldsymbol{\mu}_{-i})$ is a bounded function of $\boldsymbol{\mu}$.*

Proof. Let $v_1(\boldsymbol{\mu}) \geq v_2(\boldsymbol{\mu}) \geq \dots \geq v_n(\boldsymbol{\mu})$ be the ordered eigenvalues of Σ_μ . We can now rewrite $L(\mathbf{y}|\boldsymbol{\mu})$ as

$$L(\mathbf{y}|\boldsymbol{\mu})^2 \propto \prod_{k=1}^n \left[v_k(\boldsymbol{\mu})^{-1} (\mathbf{y}^\top \Sigma_\mu^{-1} \mathbf{y})^{-1} \right]. \quad (5.63)$$

Proposition 5.23 asserts that for any $\mathbf{y} \in \mathbb{R}^n$ that is not orthogonal to k_ν , $(\mathbf{y}^\top \Sigma_\mu^{-1} \mathbf{y})^{-1} = O(\|\boldsymbol{\mu}\|^{2\nu})$ for $\|\boldsymbol{\mu}\| \rightarrow 0$. Besides, Proposition 5.22 asserts that $\|\Sigma_\mu^{-1}\| = O(\|\boldsymbol{\mu}\|^{-2\nu})$, so $(\mathbf{y}^\top \Sigma_\mu^{-1} \mathbf{y})^{-1} = O(\|\Sigma_\mu^{-1}\|^{-1})$.

This implies that for every integer $i \in \llbracket 1, r \rrbracket$, $v_k(\boldsymbol{\mu})^{-1} (\mathbf{y}^\top \Sigma_\mu^{-1} \mathbf{y})^{-1} = O(1)$.

Clearly, $\lim_{\|\boldsymbol{\mu}\| \rightarrow 0} |\mathbf{1}^\top \mathbf{v}_1(\boldsymbol{\mu})| = \|\mathbf{1}\|$ and $\lim_{\|\boldsymbol{\mu}\| \rightarrow 0} v_1(\boldsymbol{\mu}) = n$.

The latter implies $\lim_{\|\boldsymbol{\mu}\| \rightarrow 0} v_1(\boldsymbol{\mu})^{-1} = n^{-1}$ and $v_1(\boldsymbol{\mu})^{-1} (\mathbf{y}^\top \boldsymbol{\Sigma}_\mu^{-1} \mathbf{y})^{-1} = O(\|\boldsymbol{\mu}\|^{2\nu})$. Now, for every $\boldsymbol{\mu} \in (0, +\infty)^r$, we may (thanks to the axiom of choice) choose a unit eigenvector $\mathbf{v}_1(\boldsymbol{\mu})$ corresponding to the largest eigenvalue $v_1(\boldsymbol{\mu})$ and $\mathbf{v}_2(\boldsymbol{\mu})$ corresponding to the second largest eigenvalue $v_2(\boldsymbol{\mu})$. Because $\boldsymbol{\Sigma}_\mu$ is symmetric, $\mathbf{v}_1(\boldsymbol{\mu})^\top \mathbf{v}_2(\boldsymbol{\mu}) = 0$ for all $\boldsymbol{\mu} \in (0, +\infty)^r$, so $\lim_{\|\boldsymbol{\mu}\| \rightarrow 0} \mathbf{1}^\top \mathbf{v}_2(\boldsymbol{\mu}) = 0$.

$$\begin{aligned} & v_2(\boldsymbol{\mu}) \\ &= (\mathbf{1}^\top \mathbf{v}_2(\boldsymbol{\mu}))^2 + 2\nu(\nu - 1)^{-1} \sum_{i=1}^r \mu_i^2 \left(\mathbf{X}_i^\top \mathbf{v}_2(\boldsymbol{\mu}) \right)^2 - 2\mu_i^2 (\mathbf{1}^\top \mathbf{v}_2(\boldsymbol{\mu})) \left(\mathbf{X}_i^{\circ 2\top} \mathbf{v}_2(\boldsymbol{\mu}) \right) + O(\|\boldsymbol{\mu}\|^4) \\ &\geq 2\nu(\nu - 1)^{-1} \sum_{i=1}^r \mu_i^2 (\mathbf{X}_i^\top \mathbf{v}_2(\boldsymbol{\mu}))^2 + o(\|\boldsymbol{\mu}\|^2). \end{aligned} \quad (5.64)$$

For all $\boldsymbol{\mu} \in (0, +\infty)^r$, let $i(\boldsymbol{\mu})$ be the smallest integer $i \in \llbracket 1, r \rrbracket$ such that $\mu_i = \max_{j=1}^r \mu_j$. Now for every integer $i \in \llbracket 1, r \rrbracket$ let $\mathbf{w}_i(\boldsymbol{\mu})$ be the unit vector that belongs to the space spanned by $\mathbf{v}_1(\boldsymbol{\mu})$ and \mathbf{X}_i that verifies $\mathbf{v}_1(\boldsymbol{\mu})^\top \mathbf{w}_i(\boldsymbol{\mu}) = 0$ and $\mathbf{X}_i^\top \mathbf{w}_i(\boldsymbol{\mu}) > 0$.

$$\mathbf{w}_{i(\boldsymbol{\mu})}(\boldsymbol{\mu})^\top \boldsymbol{\Sigma}_\mu \mathbf{w}_{i(\boldsymbol{\mu})}(\boldsymbol{\mu}) \geq 2\nu(\nu - 1)^{-1} r^{-1} \|\boldsymbol{\mu}\|^2 (\mathbf{X}_{i(\boldsymbol{\mu})}^\top \mathbf{w}_{i(\boldsymbol{\mu})}(\boldsymbol{\mu}))^2 + o(\|\boldsymbol{\mu}\|^2). \quad (5.65)$$

As $\lim_{\|\boldsymbol{\mu}\| \rightarrow 0} |\mathbf{1}^\top \mathbf{v}_1(\boldsymbol{\mu})| = \|\mathbf{1}\|$, $\liminf_{\|\boldsymbol{\mu}\| \rightarrow 0} \mathbf{X}_{i(\boldsymbol{\mu})}^\top \mathbf{w}_{i(\boldsymbol{\mu})}(\boldsymbol{\mu}) \geq \min_{i=1}^r \lim_{\|\boldsymbol{\mu}\| \rightarrow 0} \mathbf{X}_i^\top \mathbf{w}_i(\boldsymbol{\mu}) > 0$, so there exists a constant $c_2 > 0$ such that when $\|\boldsymbol{\mu}\|$ is small enough

$$\mathbf{w}_{i(\boldsymbol{\mu})}(\boldsymbol{\mu})^\top \boldsymbol{\Sigma}_\mu \mathbf{w}_{i(\boldsymbol{\mu})}(\boldsymbol{\mu}) \geq c_2 \|\boldsymbol{\mu}\|^2. \quad (5.66)$$

Recall $v_2(\boldsymbol{\mu}) = \max\{\boldsymbol{\xi}^\top \boldsymbol{\Sigma}_\mu \boldsymbol{\xi} \mid \boldsymbol{\xi} \in S^{n-1} \text{ and } \boldsymbol{\xi}^\top \mathbf{v}_1(\boldsymbol{\mu})\}$, so a fortiori $v_2(\boldsymbol{\mu}) \geq c_2 \|\boldsymbol{\mu}\|^2$.

This implies $v_2(\boldsymbol{\mu})^{-1} = O(\|\boldsymbol{\mu}\|^{-2})$ and therefore $v_2(\boldsymbol{\mu})^{-1} (\mathbf{y}^\top \boldsymbol{\Sigma}_\mu^{-1} \mathbf{y})^{-1} = O(\|\boldsymbol{\mu}\|^{2(\nu-1)})$.

Finally, $L(\mathbf{y}|\boldsymbol{\mu}) = O(\|\boldsymbol{\mu}\|^\nu) O(\|\boldsymbol{\mu}\|^{\nu-1}) = O(\|\boldsymbol{\mu}\|^{2\nu-1})$.

Given that $f_i(\mu_i|\boldsymbol{\mu}_{-i}) = O(\|\boldsymbol{\mu}\|^{1-2\nu})$, $L(\mathbf{y}|\boldsymbol{\mu})f_i(\mu_i|\boldsymbol{\mu}_{-i})$ is bounded when $\|\boldsymbol{\mu}\| \rightarrow 0$. \square

Proposition 5.25. *Assume $\nu \in (0, 1)$. For every $\mathbf{y} \in \mathbb{R}^n$ that is not collinear to $\mathbf{1}$, when $\|\boldsymbol{\mu}\| \rightarrow 0$, $L(\mathbf{y}|\boldsymbol{\mu})f_i(\mu_i|\boldsymbol{\mu}_{-i}) = O(\mu_i^{-1+\nu})$.*

Proof. This proof is similar to the previous one, so we use the same notations. $v_1(\boldsymbol{\mu})^{-1} = O(1)$, so $v_1(\boldsymbol{\mu})^{-1} (\mathbf{y}^\top \boldsymbol{\Sigma}_\mu^{-1} \mathbf{y})^{-1} = O(\|\boldsymbol{\mu}\|^{2\nu})$, which yields that $L(\mathbf{y}|\boldsymbol{\mu}) = O(\|\boldsymbol{\mu}\|^\nu)$. This implies that $L(\mathbf{y}|\boldsymbol{\mu}) \|\boldsymbol{\Sigma}_\mu^{-1}\| = O(\|\boldsymbol{\mu}\|^{-\nu}) = O(\mu_i^{-\nu})$.

By Lemma 5.7, $\left\| \frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_\mu \right\| = O(\mu_i^{-1+2\nu})$. Putting all this together, $L(\mathbf{y}|\boldsymbol{\mu})f_i(\mu_i|\boldsymbol{\mu}_{-i}) = O(\mu_i^{-1+\nu})$. \square

Proposition 5.26. *For Matérn anisotropic geometric or tensorized kernels with smoothness parameter $\nu \in (0, 1) \cup (1, 2) \cup (2, 3)$, if $\mathbf{y} \in \mathbb{R}^n$ is not orthogonal to k_ν , then the conditional posterior distribution $\pi_i(\mu_i|\mathbf{y}, \boldsymbol{\mu}_{-i})$, seen as a function of $\boldsymbol{\mu}$, is continuous over $\{\boldsymbol{\mu} \in [0, +\infty)^r : \mu_i \neq 0\}$.*

Moreover,

$$\forall \mu_i > 0, \quad \pi_i(\mu_i|\mathbf{y}, \boldsymbol{\mu}_{-i} = \mathbf{0}_{r-1}) = \frac{L(\mathbf{y}|\mu_i, \boldsymbol{\mu}_{-i} = \mathbf{0}_{r-1})f_i(\mu_i|\boldsymbol{\mu}_{-i} = \mathbf{0}_{r-1})}{\int_0^\infty L(\mathbf{y}|\mu_i = t, \boldsymbol{\mu}_{-i} = \mathbf{0}_{r-1})f_i(\mu_i = t|\boldsymbol{\mu}_{-i} = \mathbf{0}_{r-1})dt} > 0.$$

Proof. Given Proposition 5.8 and the fact that $\forall \mathbf{y} \in \mathbb{R}^n$, $\forall i \in \llbracket 1, r \rrbracket$ and $\forall \mu_i \in (0, +\infty)$, as $\|\boldsymbol{\mu}_{-i}\| \rightarrow 0$, $L(\mathbf{y}|\boldsymbol{\mu})f_i(\mu_i|\boldsymbol{\mu}_{-i})$ converges pointwise to $L(\mathbf{y}|\mu_i, \boldsymbol{\mu}_{-1} = \mathbf{0}_{r-1})f_i(\mu_i|\boldsymbol{\mu}_{-i} = \mathbf{0}_{r-1})$, we only need to show that

$$\begin{aligned} \int_0^\infty L(\mathbf{y}|\mu_i = t, \boldsymbol{\mu}_{-i})f_i(\mu_i = t|\boldsymbol{\mu}_{-i})dt &\xrightarrow{\|\boldsymbol{\mu}_{-i}\| \rightarrow 0} \\ \int_0^\infty L(\mathbf{y}|\mu_i = t, \boldsymbol{\mu}_{-i} = \mathbf{0}_{r-1})f_i(\mu_i = t|\boldsymbol{\mu}_{-i} = \mathbf{0}_{r-1})dt &< +\infty. \end{aligned} \quad (5.67)$$

Lemma 5.7 implies that there exists $M_i > 0$ such that

$$\begin{aligned} L(\mathbf{y}|\boldsymbol{\mu})f_i(\mu_i|\boldsymbol{\mu}_{-i}) &= L(\mathbf{y}|\boldsymbol{\mu}) \|\boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{-1}\| \|\boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{-1}\|^{-1} f_i(\mu_i|\boldsymbol{\mu}_{-i}) \\ &\leq L(\mathbf{y}|\boldsymbol{\mu}) \|\boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{-1}\| M_i \mu_i^{-2}. \end{aligned} \quad (5.68)$$

Lemma 5.10 then ensures $\|\boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{-1} - \mathbf{I}_n\| \rightarrow 0$ as $\mu_i \rightarrow +\infty$, so the right member of the inequality is integrable in the neighborhood of $+\infty$. Let us now focus on the neighborhood of 0.

If $\nu > 1$, Proposition 5.24 asserts that $L(\mathbf{y}|\boldsymbol{\mu})f_i(\mu_i|\boldsymbol{\mu}_{-i})$ is bounded in the neighborhood of 0.

If $\nu < 1$, Proposition 5.25 asserts that $L(\mathbf{y}|\boldsymbol{\mu})f_i(\mu_i|\boldsymbol{\mu}_{-i})\mu_i^{1-\nu}$ is bounded in the neighborhood of 0.

Therefore, there exists a function independent of $\boldsymbol{\mu}_{-i}$ that is both greater than the product $L(\mathbf{y}|\boldsymbol{\mu})f_i(\mu_i|\boldsymbol{\mu}_{-i})$ and integrable over $\mu_i \in (0, +\infty)$, so the dominated convergence theorem is applicable. □

Lower bound for conditional reference posterior densities

The following Lemma provides the key to proving Theorem 5.2:

Lemma 5.27. *In a Simple Kriging model with the characteristics described above, there exists a hyperplane \mathcal{H} of \mathbb{R}^n such that, for any $\mathbf{y} \in \mathbb{R}^n \setminus \mathcal{H}$ and any $i \in \llbracket 1, r \rrbracket$, there exists a measurable function $m_{i,\mathbf{y}} : (0, +\infty) \rightarrow (0, +\infty)$ such that, for all $\boldsymbol{\mu}_{-i} \in (0, +\infty)^{r-1}$, the conditional reference posterior density verifies:*

$$\pi_i(\mu_i|\mathbf{y}, \boldsymbol{\mu}_{-i}) \geq m_{i,\mathbf{y}}(\mu_i) > 0. \quad (5.69)$$

Proof. This proof consists in combining Proposition 5.19 and Proposition 5.26, which respectively deal with large and small values of $\|\boldsymbol{\mu}_{-i}\|$.

Proposition 5.19 implies that for any $\mathbf{y} \in \mathbb{R}^n \setminus \{0\}^n$, for any $i \in \llbracket 1, r \rrbracket$ and any $\mu_i \in (0, +\infty)$, there exists a compact neighborhood $N_i(\mu_i)$ of $\mathbf{0}_{r-1}$ within $[0, +\infty)^r$ such that

$$\inf\{\pi(\mu_i|\mathbf{y}, \boldsymbol{\mu}_{-i}) : \boldsymbol{\mu}_{-i} \in [0, +\infty)^{r-1} \setminus N_i(\mu_i)\} > 0. \quad (5.70)$$

The vector space $k_\nu \subset \mathbb{R}^n$ has dimension greater or equal to

- (a) $n - 1$ if $\nu \in (0, 1)$;
- (b) $n - (r + 1)$ if $\nu \in (1, 2)$;

(c) $n - (r + 1)(r + 2)/2$ if $\nu \in (2, 3)$.

For all Simple Kriging models tackled by this lemma, the dimension of k_ν is therefore greater or equal to 1. Its orthogonal complement k_ν^\perp is then included within a hyperplane \mathcal{H} of \mathbb{R}^n . Assuming $\mathbf{y} \in \mathbb{R}^n \setminus \mathcal{H}$, Proposition 5.26 ensures that $\pi(\mu_i | \mathbf{y}, \boldsymbol{\mu}_{-i})$ is a continuous and positive function of $\boldsymbol{\mu}$ on $\{\boldsymbol{\mu} \in [0, +\infty)^r : \mu_i \neq 0\}$. In particular, this implies that for any $\mu_i \in (0, +\infty)$ and any compact neighborhood $N_i(\mu_i)$ of $\mathbf{0}_{r-1}$ within $[0, +\infty)^r$,

$$\inf\{\pi(\mu_i | \mathbf{y}, \boldsymbol{\mu}_{-i}) : \boldsymbol{\mu}_{-i} \in N_i(\mu_i)\} > 0. \quad (5.71)$$

Putting this together, if $\mathbf{y} \in \mathbb{R}^n \setminus \mathcal{H}$, for any $i \in \llbracket 1, r \rrbracket$ and any $\mu_i \in (0, +\infty)$,

$$\inf\{\pi(\mu_i | \mathbf{y}, \boldsymbol{\mu}_{-i}) : \boldsymbol{\mu}_{-i} \in [0, +\infty)^{r-1}\} > 0. \quad (5.72)$$

The mapping $m_{i,\mathbf{y}} : \mu_i \mapsto \inf\{\pi(\mu_i | \mathbf{y}, \boldsymbol{\mu}_{-i}) : \boldsymbol{\mu}_{-i} \in [0, +\infty)^{r-1}\}$, which is measurable on $(0, +\infty)$ is therefore also positive on $(0, +\infty)$. □

Proof of Theorem 5.2. Lemma 5.27 implies that $\forall \boldsymbol{\mu}^{(0)} \in (0, +\infty)^r$,

$$P_{\mathbf{y}}(\boldsymbol{\mu}^{(0)}, d\boldsymbol{\mu}) \geq \frac{1}{r} \sum_{i=1}^r m_{i,\mathbf{y}}(\mu_i) d\mu_i \delta_{\boldsymbol{\mu}_{-i}^{(0)}}(d\boldsymbol{\mu}_{-i}),$$

and thus $\forall \boldsymbol{\mu}^{(0)} \in (0, +\infty)^r$, $\forall n \geq r$,

$$P_{\mathbf{y}}^n(\boldsymbol{\mu}^{(0)}, d\boldsymbol{\mu}) \geq \frac{1}{r^n} \prod_{i=1}^r m_{i,\mathbf{y}}(\mu_i) d\mu_i.$$

Defining $f_{\mathbf{y}}(\boldsymbol{\mu}) := r^{-r} \prod_{i=1}^r m_{i,\mathbf{y}}(\mu_i)$, $f_{\mathbf{y}}$ is a measurable positive function. Therefore $f_{\mathbf{y}}$ is the density with respect to the Lebesgue measure of a positive measure with mass $\epsilon_{\mathbf{y}} > 0$. So $\epsilon_{\mathbf{y}}^{-1} f_{\mathbf{y}}$ is a probability density with respect to the Lebesgue measure and the Markov kernel $P_{\mathbf{y}}$ thus satisfies the uniform $(n, \epsilon_{\mathbf{y}})$ Doeblin condition:

$$\forall \boldsymbol{\mu}^{(0)} \in (0, +\infty)^r \quad P_{\mathbf{y}}^n(\boldsymbol{\mu}^{(0)}, d\boldsymbol{\mu}) \geq \epsilon_{\mathbf{y}} \left(\frac{1}{\epsilon_{\mathbf{y}}} f_{\mathbf{y}}(\boldsymbol{\mu}) \right) d\boldsymbol{\mu}. \quad (5.73)$$

This implies that $P_{\mathbf{y}}$ is uniformly ergodic: it has a unique invariant probability distribution $\pi_G(\cdot | \mathbf{y})$ and $\lim_{n \rightarrow \infty} \sup_{\boldsymbol{\mu}^{(0)} \in (0, +\infty)^r} \|P_{\mathbf{y}}^n(\boldsymbol{\mu}^{(0)}, \cdot) - \pi_G(\cdot | \mathbf{y})\|_{TV} = 0$, where $\|\cdot\|_{TV}$ is the total variation norm. By definition, $\pi_G(\cdot | \mathbf{y})$ is the Gibbs compromise between the incompatible posterior conditionals $\pi_i(\mu_i | \mathbf{y}, \boldsymbol{\mu}_{-i})$. □

Chapter 6

A Comprehensive Bayesian Treatment of the Universal Kriging model with Matérn correlation kernels

This chapter mostly adheres to the article Muré [2018].

Abstract

The Gibbs reference posterior distribution provides an objective full-Bayesian solution to the problem of prediction of a stationary Gaussian process with Matérn anisotropic kernel. A full-Bayesian approach is possible, because the posterior distribution is expressed as the invariant distribution of a uniformly ergodic Markovian kernel for which we give an explicit expression. In this chapter, we show that it is appropriate for the Universal Kriging framework, that is when an unknown function is added to the stationary Gaussian process. We give sufficient conditions for the existence and propriety of the Gibbs reference posterior that apply to a wide variety of practical cases and illustrate the method with several examples. Finally, simulations of Gaussian processes suggest that the Gibbs reference posterior has good frequentist properties in terms of coverage of prediction intervals.

Résumé

Le posterior de référence de Gibbs fournit une solution pleinement bayésienne au problème de la prédiction des valeurs prises par un processus gaussien stationnaire avec noyau de Matérn anisotrope. Une approche pleinement bayésienne est possible parce que la loi *a posteriori* est exprimée comme loi invariante d'un noyau markovien uniformément ergodique dont nous donnons une expression explicite. Dans ce chapitre, nous montrons que cette approche est indiquée dans le cadre du krigeage universel, c'est-à-dire quand une fonction inconnue est ajoutée à un processus gaussien stationnaire. Nous donnons des conditions suffisantes à l'existence et la propriété du posterior de référence de Gibbs qui s'appliquent à une large collection de cas pratiques et illustrons la méthode par plusieurs exemples. Enfin, des simulation des processus gaussien suggèrent que le posterior de référence de Gibbs a de bonnes propriétés fréquentistes en termes de couverture et d'intervalles prédictifs.

6.1 Introduction

In Simple Kriging models, the Gaussian Process is assumed to have zero mean and be stationary on the domain \mathcal{D} , so its distribution can be characterized by a positive variance parameter σ^2 and by an autocorrelation function K . The Universal Kriging framework adds another parameter: a mean function f . If f is known, then subtracting it from the process returns us to the Simple Kriging framework. Allowing for an unknown mean function provides greater flexibility in the modeling by enabling some degree of non-stationarity [Santner et al., 2003, section 2.3.2].

In practice, the mean function f is assumed to belong to a p -dimensional ($p \in \mathbb{N}$) vector space \mathcal{F}_p , which is specified by means of a basis (f_1, \dots, f_p) . Being a linear combination of f_1, \dots, f_p , the mean function f is then encoded by the vector of linear coefficients $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$: $f = \beta_1 f_1 + \dots + \beta_p f_p$.

Therefore, what separates the Universal Kriging framework from its Simple counterpart is the addition of the p -dimensional parameter $\boldsymbol{\beta}$.

Assuming the autocorrelation function to be characterized by a vector of correlation lengths $\boldsymbol{\theta}$, we are faced with the inference problem of estimating $(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\theta})$. In the previous chapter, an objective posterior distribution on $(\sigma^2, \boldsymbol{\theta})$ was proposed in the context of Simple Kriging. In this chapter we address the more general framework of Universal Kriging in order to obtain a distribution on $(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\theta})$. The developments in both chapters are based on Bernardo's reference prior theory. The idea to use this theory in the context of Kriging first appeared in Berger et al. [2001], and then was successively extended by Paulo [2005], Kazianka and Pilz [2012], Ren et al. [2012], Ren et al. [2013] and Gu [2016].

To use it we first need to order the parameters [Bernardo, 2005]. Because our main goal is to maximize the predictive capacity of the model, we are unable to outright say which parameter we care about most. However, a few common sense observations help: first, in order to profit from the work done in the Simple Kriging case, we separate $\boldsymbol{\beta}$, which refers to the mean function, from $(\sigma^2, \boldsymbol{\theta})$, which yields the covariance structure. Within the latter, $\boldsymbol{\theta}$ should have the priority over σ^2 , because while σ^2 can very easily be accurately estimated once $\boldsymbol{\theta}$ is known, the reverse is not true. The same consideration will make us prioritize $(\sigma^2, \boldsymbol{\theta})$ over $\boldsymbol{\beta}$, because while knowing $\boldsymbol{\beta}$ reduces the problem to the Simple Kriging case, knowing $(\sigma^2, \boldsymbol{\theta})$ reduces it to a much simpler regression problem.

In Section 6.2 we derive the reference posterior distribution on $(\boldsymbol{\beta}, \sigma^2)$ and the corresponding predictive distribution at unobserved points, both conditional to the observed data and the correlation parameter $\boldsymbol{\theta}$.

In Section 6.3, we derive analytical formulas for the reference prior on $(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\theta})$ in the case where $\boldsymbol{\theta}$ is a one-dimensional parameter.

In Section 6.4, we prove the main result of the chapter: in the context of a Matérn anisotropic correlation kernel [Matérn, 1986, Handcock and Stein, 1993] – see Appendix 6.A for precise definitions – under a few conditions, the Gibbs reference posterior on a multidimensional $\boldsymbol{\theta}$ exists. Combined with the “partial” reference posterior distribution on $(\boldsymbol{\beta}, \sigma^2)$ conditional to $\boldsymbol{\theta}$, it provides a proper objective posterior distribution on all parameters $(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\theta})$ given the observed data. It is significant that this proper objective posterior distribution is well

defined for Matérn anisotropic correlation kernels, because this class of correlation kernels has remarkable properties (see Stein [1999] or Bachoc [2013a, chapter 2]). Notably, it allows the user to specify the smoothness of the realizations of the Gaussian Process.

In Section 6.5, we evaluate the predictive performance of the Universal Kriging model with the Gibbs reference posterior distribution both in the context of a well-specified model and when emulating deterministic functions. We compare the full-Bayesian approach relying on the Gibbs reference posterior with plug-in approaches, where the parameters are assumed to be equal to either the Maximum Likelihood Estimator (MLE) or the Maximum A Posteriori (MAP) estimator.

6.2 Analytical treatment of the location-scale parameters

Suppose our design set contains n observation points. n must be greater than p , otherwise the model is not identifiable. Let \mathbf{H} be the $n \times p$ matrix whose columns contain the values of the p basis functions at the n observation points. Let us assume that the rank of \mathbf{H} is p , because if it were not, the model would also not be identifiable.

Let \mathbf{y} be the vector of the n observations. Then \mathbf{y} is a Gaussian vector and its distribution is

$$\mathbf{y} | \boldsymbol{\beta}, \sigma^2, \boldsymbol{\theta} \sim \mathcal{N}(\mathbf{H}\boldsymbol{\beta}, \sigma^2 \boldsymbol{\Sigma}_{\boldsymbol{\theta}}), \quad (6.1)$$

where $\boldsymbol{\Sigma}_{\boldsymbol{\theta}}$ is a correlation matrix that only depends on the design set and on the vector of correlation lengths $\boldsymbol{\theta}$.

In terms of likelihood, we have

$$L(\mathbf{y} | \boldsymbol{\beta}, \sigma^2, \boldsymbol{\theta}) = \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{n}{2}} |\boldsymbol{\Sigma}_{\boldsymbol{\theta}}|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{H}\boldsymbol{\beta})^\top \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} (\mathbf{y} - \mathbf{H}\boldsymbol{\beta}) \right\}. \quad (6.2)$$

The aim of this section is to get the parameters $\boldsymbol{\beta}$ and σ^2 out of the way in order to focus on the more interesting parameter $\boldsymbol{\theta}$. For now, assume that $\boldsymbol{\theta}$ is known, which is to say that the correlation function is completely known.

Reference prior and integrated likelihood when $\boldsymbol{\theta}$ is known.

Clearly, $\boldsymbol{\beta}$ is a location parameter and $\sigma := \sqrt{\sigma^2}$ is a scale parameter for this model. Therefore, the joint reference prior is $\pi(\boldsymbol{\beta}, \sigma^2 | \boldsymbol{\theta}) \propto 1/\sigma^2$ regardless of the order of the parameters $(\boldsymbol{\beta}, \sigma^2)$.

We now derive the posterior distributions $\pi(\boldsymbol{\beta} | \mathbf{y}, \sigma^2, \boldsymbol{\theta})$ and $\pi(\sigma^2 | \mathbf{y}, \boldsymbol{\theta})$ as well as the integrated likelihoods $L^0(\mathbf{y} | \sigma^2, \boldsymbol{\theta}) := \int L(\mathbf{y} | \boldsymbol{\beta}, \sigma^2, \boldsymbol{\theta}) d\boldsymbol{\beta}$ and $L^1(\mathbf{y} | \boldsymbol{\theta}) := \iint L(\mathbf{y} | \boldsymbol{\beta}, \sigma^2, \boldsymbol{\theta}) / \sigma^2 d\boldsymbol{\beta} d\sigma^2$.

Gaussian theory makes it convenient to split \mathbf{y} into two components: one that belongs to the subspace of \mathbb{R}^n spanned by \mathbf{H} , and one that is orthogonal to the subspace spanned by \mathbf{H} . In order not to have to deal with degenerate Gaussian vectors, we define an $n \times p$ matrix \mathbf{P} with full rank which spans the same subspace as \mathbf{H} (Actually, for the time being, we may as well set $\mathbf{P} = \mathbf{H}$.) and an $n \times (n - p)$ matrix \mathbf{W} with full rank which spans its orthogonal space. Thus $\mathbf{W}^\top \mathbf{H} = \mathbf{W}^\top \mathbf{P} = \mathbf{0}_{n-p,p}$ and

$$\mathbf{W}^\top \mathbf{y} | \sigma^2, \boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}_{n-p}, \sigma^2 \mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W}); \quad (6.3)$$

$$\begin{aligned} \mathbf{P}^\top \mathbf{y} | \beta, \sigma^2, \boldsymbol{\theta}, \mathbf{W}^\top \mathbf{y} &\sim \mathcal{N}(\mathbf{P}^\top \mathbf{H} \beta + \mathbf{P}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} (\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{y}, \\ &\sigma^2 \mathbf{P}^\top \boldsymbol{\Sigma}_\theta \mathbf{P} - \sigma^2 \mathbf{P}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} (\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W})^{-1} \mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{P}). \end{aligned} \quad (6.4)$$

β having flat prior density, $\mathbf{P}^\top \mathbf{y} - \mathbf{P}^\top \mathbf{H} \beta$ has the same distribution whether β , σ^2 and $\boldsymbol{\theta}$ or whether $\mathbf{P}^\top \mathbf{y}$, σ^2 and $\boldsymbol{\theta}$ are known. Therefore, the posterior distribution of $\mathbf{P}^\top \mathbf{H} \beta$ if σ^2 and $\boldsymbol{\theta}$ are known is:

$$\begin{aligned} \mathbf{P}^\top \mathbf{H} \beta | \sigma^2, \boldsymbol{\theta}, \mathbf{W}^\top \mathbf{y}, \mathbf{P}^\top \mathbf{y} &\sim \mathcal{N}(\mathbf{P}^\top \mathbf{y} - \mathbf{P}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} (\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{y}, \\ &\sigma^2 \mathbf{P}^\top \boldsymbol{\Sigma}_\theta \mathbf{P} - \sigma^2 \mathbf{P}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} (\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W})^{-1} \mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{P}). \end{aligned} \quad (6.5)$$

From there, we get the posterior distribution of β if σ^2 and $\boldsymbol{\theta}$ are known:

$$\begin{aligned} \beta | \sigma^2, \boldsymbol{\theta}, \mathbf{y} &\sim \mathcal{N}((\mathbf{P}^\top \mathbf{H})^{-1} \mathbf{P}^\top \mathbf{y} - (\mathbf{P}^\top \mathbf{H})^{-1} \mathbf{P}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} (\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{y}, \\ &\sigma^2 ((\mathbf{P}^\top \mathbf{H})^{-1} \mathbf{P}^\top \boldsymbol{\Sigma}_\theta \mathbf{P} (\mathbf{H}^\top \mathbf{P})^{-1} \\ &- (\mathbf{P}^\top \mathbf{H})^{-1} \mathbf{P}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} (\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W})^{-1} \mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{P} (\mathbf{H}^\top \mathbf{P})^{-1}) \end{aligned} \quad (6.6)$$

Moreover, (6.4) implies that the integrated likelihood of $\mathbf{P}^\top \mathbf{y}$, i.e. its likelihood averaged over the Lebesgue measure (the prior distribution on β), is $|\mathbf{P}^\top \mathbf{H}|^{-1}$, where $|\cdot|$ denotes the absolute value of the determinant.

$$\mathbf{P}^\top \mathbf{y} | \sigma^2, \boldsymbol{\theta}, \mathbf{W}^\top \mathbf{y} \sim \text{Improper "uniform" distribution on } \mathbb{R}^p. \quad (6.7)$$

This means that if β is unknown, then $\mathbf{P}^\top \mathbf{y}$ can yield no information about σ^2 and $\boldsymbol{\theta}$. When β is unknown, all information about σ^2 and $\boldsymbol{\theta}$ is carried by $\mathbf{W}^\top \mathbf{y}$, because as is shown by (6.3), the predictive distribution on $\mathbf{W}^\top \mathbf{y}$ knowing σ^2 and $\boldsymbol{\theta}$ does not depend on β .

A straightforward calculation yields that the posterior distribution of σ^2 is Inverse-Gamma:

$$\sigma^2 | \boldsymbol{\theta}, \mathbf{y} \sim \mathcal{IG}(\text{shape} = (n - p)/2, \text{rate} = \mathbf{y}^\top \mathbf{W} (\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{y}/2). \quad (6.8)$$

The posterior distribution on σ^2 (knowing $\boldsymbol{\theta}$) does not take into account $\mathbf{P}^\top \mathbf{y}$, because all information contained in $\mathbf{P}^\top \mathbf{y}$ is given in the posterior distribution of β conditional to σ^2 and $\boldsymbol{\theta}$.

We conclude this subsection with the formulas for the likelihoods with the parameters β and σ^2 successively integrated out.

$$\begin{aligned}
L^0(\mathbf{y}|\sigma^2, \boldsymbol{\theta}) &= \int L(\mathbf{y}|\boldsymbol{\beta}, \sigma^2, \boldsymbol{\theta}) d\boldsymbol{\beta} \\
&= |\mathbf{P}^\top \mathbf{H}|^{-1} \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{n-p}{2}} |\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W}|^{-\frac{1}{2}} \exp \left\{ -\frac{\mathbf{y}^\top \mathbf{W} (\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{y}}{2\sigma^2} \right\}; \tag{6.9}
\end{aligned}$$

$$\begin{aligned}
L^1(\mathbf{y}|\boldsymbol{\theta}) &= \int L^0(\mathbf{y}|\sigma^2, \boldsymbol{\theta}) / \sigma^2 d\sigma^2 \\
&= |\mathbf{P}^\top \mathbf{H}|^{-1} \left(\frac{2\pi^{\frac{n-p}{2}}}{\Gamma(\frac{n-p}{2})} \right)^{-1} |\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W}|^{-\frac{1}{2}} \left(\mathbf{y}^\top \mathbf{W} (\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{y} \right)^{-\frac{n-p}{2}}. \tag{6.10}
\end{aligned}$$

Posterior predictive distribution when $\boldsymbol{\theta}$ is known.

Following Santner et al. [2003] (Theorem 4.1.2., case (4)), we derive conditionally to $\boldsymbol{\theta}$ the posterior predictive distribution of the values taken by the process at unobserved points.

In order to simplify notations in this subsection, all the distributions we consider are, until further notice, conditional to σ^2 and $\boldsymbol{\theta}$ even with no explicit mention. Equation (6.1) can be usefully restated in the following way:

$$\left(\begin{array}{c} \mathbf{P}^\top \mathbf{y} - \mathbf{P}^\top \mathbf{H}\boldsymbol{\beta} \\ \mathbf{W}^\top \mathbf{y} \end{array} \right) \middle| \mathbf{P}^\top \mathbf{H}\boldsymbol{\beta} \sim \mathcal{N} \left(\mathbf{0}_n, \sigma^2 \begin{pmatrix} \mathbf{P}^\top \\ \mathbf{W}^\top \end{pmatrix} \boldsymbol{\Sigma}_\theta \begin{pmatrix} \mathbf{P} & \mathbf{W} \end{pmatrix} \right). \tag{6.11}$$

Because the prior distribution on $\boldsymbol{\beta}$ is flat, $\left(\begin{array}{c} \mathbf{P}^\top \mathbf{y} - \mathbf{P}^\top \mathbf{H}\boldsymbol{\beta} \\ \mathbf{W}^\top \mathbf{y} \end{array} \right)$ and its opposite have the same distribution when conditional respectively to $\mathbf{P}^\top \mathbf{H}\boldsymbol{\beta}$ and $\mathbf{P}^\top \mathbf{y}$.

$$\left(\begin{array}{c} \mathbf{P}^\top \mathbf{H}\boldsymbol{\beta} - \mathbf{P}^\top \mathbf{y} \\ -\mathbf{W}^\top \mathbf{y} \end{array} \right) \middle| \mathbf{P}^\top \mathbf{y} \sim \mathcal{N} \left(\mathbf{0}_n, \sigma^2 \begin{pmatrix} \mathbf{P}^\top \\ \mathbf{W}^\top \end{pmatrix} \boldsymbol{\Sigma}_\theta \begin{pmatrix} \mathbf{P} & \mathbf{W} \end{pmatrix} \right) \tag{6.12}$$

Let \mathbf{y}_0 be the values of the Gaussian Process at the n_0 unobserved points. We denote by $\mathbf{H}_{0,0}$ the $n_0 \times p$ matrix whose columns contain the values of the p basis functions at the unobserved points, by $\boldsymbol{\Sigma}_{\theta,0,0}$ the $n_0 \times n_0$ correlation matrix of \mathbf{y}_0 , by $\boldsymbol{\Sigma}_{\theta,0,\cdot}$ the $n_0 \times n$ correlation matrix between \mathbf{y}_0 and \mathbf{y} and by $\boldsymbol{\Sigma}_{\theta,\cdot,0}$ its transpose. It is also convenient to define the $n_0 \times n$ matrix $\mathbf{H}_{0,\cdot} = \mathbf{H}_{0,0} (\mathbf{P}^\top \mathbf{H})^{-1}$ and its transpose $\mathbf{H}_{\cdot,0}$. With these notations, the distribution of \mathbf{y}_0 when \mathbf{y} and $\boldsymbol{\beta}$ are known is

$$\mathbf{y}_0 | \boldsymbol{\beta}, \mathbf{y} \sim \mathcal{N} (\mathbf{H}_{0,0}\boldsymbol{\beta} + \boldsymbol{\Sigma}_{\theta,0,\cdot} \boldsymbol{\Sigma}_\theta^{-1} (\mathbf{y} - \mathbf{H}\boldsymbol{\beta}), \boldsymbol{\Sigma}_{\theta,0,0} - \boldsymbol{\Sigma}_{\theta,0,\cdot} \boldsymbol{\Sigma}_\theta^{-1} \boldsymbol{\Sigma}_{\theta,\cdot,0}) \tag{6.13}$$

Now, the distribution of \mathbf{y}_0 when both \mathbf{y} and $\boldsymbol{\beta}$ are known, together with the distribution of $\left(\begin{array}{c} \mathbf{P}^\top \mathbf{H}\boldsymbol{\beta} - \mathbf{P}^\top \mathbf{y} \\ -\mathbf{W}^\top \mathbf{y} \end{array} \right)$ when $\mathbf{P}^\top \mathbf{y}$ is known, jointly define some probability distribution on the vector $\left(\begin{array}{c} \mathbf{y}_0 \\ \mathbf{P}^\top \mathbf{H}\boldsymbol{\beta} - \mathbf{P}^\top \mathbf{y} \\ -\mathbf{W}^\top \mathbf{y} \end{array} \right)$.

This distribution is given in the following proposition. In order to give it a concise expression, it is convenient to require that $\mathbf{P}\mathbf{P}^\top + \mathbf{W}\mathbf{W}^\top = \mathbf{I}_n$, which simply means that the columns of \mathbf{P} and \mathbf{W} form an orthonormal basis of $\langle \mathbf{H} \rangle$ and its orthogonal space respectively.

Proposition 6.1. *Assume that $\mathbf{P}\mathbf{P}^\top + \mathbf{W}\mathbf{W}^\top = \mathbf{I}_n$. Then the probability distribution on the vector of \mathbb{R}^{n_0+n} $(\mathbf{y}_0, \mathbf{P}^\top \mathbf{H}\boldsymbol{\beta} - \mathbf{P}^\top \mathbf{y}, -\mathbf{W}^\top \mathbf{y})^\top$ conditional to $\mathbf{P}^\top \mathbf{y}$ is the following multivariate normal distribution:*

$$\mathcal{N} \left(\begin{pmatrix} \mathbf{E}_0 \\ \mathbf{0}_n \end{pmatrix}, \sigma^2 \begin{pmatrix} \mathbf{S}_{\theta,0,0} & \mathbf{S}_{\theta,0,\cdot} \\ \mathbf{S}_{\theta,\cdot,0} & \begin{pmatrix} \mathbf{P}^\top \\ \mathbf{W}^\top \end{pmatrix} \boldsymbol{\Sigma}_\theta \begin{pmatrix} \mathbf{P} & \mathbf{W} \end{pmatrix} \end{pmatrix} \right). \quad (6.14)$$

We use the following notations:

$$\begin{aligned} \mathbf{E}_0 &:= \mathbf{H}_{0,\cdot} \mathbf{P}^\top \mathbf{y} \\ \mathbf{S}_{\theta,0,0} &:= \boldsymbol{\Sigma}_{\theta,0,0} + \mathbf{H}_{0,\cdot} \mathbf{P}^\top \boldsymbol{\Sigma}_\theta \mathbf{P} \mathbf{H}_{\cdot,0} - \mathbf{H}_{0,\cdot} \mathbf{P}^\top \boldsymbol{\Sigma}_{\theta,\cdot,0} - \boldsymbol{\Sigma}_{\theta,0,\cdot} \mathbf{P} \mathbf{H}_{\cdot,0} \\ \mathbf{S}_{\theta,0,\cdot} \begin{pmatrix} \mathbf{P}^\top \\ \mathbf{W}^\top \end{pmatrix} &:= \mathbf{H}_{0,\cdot} \mathbf{P}^\top \boldsymbol{\Sigma}_\theta - \boldsymbol{\Sigma}_{\theta,0,\cdot} \\ \mathbf{S}_{\theta,\cdot,0} &:= \mathbf{S}_{\theta,0,\cdot}^\top \end{aligned}$$

Proof. First, notice that the mean vector of the Normal distribution given by Equation 6.13 can be rewritten as

$$(\mathbf{H}_{0,\cdot} - \boldsymbol{\Sigma}_{\theta,0,\cdot} \boldsymbol{\Sigma}_\theta^{-1} \mathbf{P}) \mathbf{P}^\top \mathbf{H}\boldsymbol{\beta} + \boldsymbol{\Sigma}_{\theta,0,\cdot} \boldsymbol{\Sigma}_\theta^{-1} \mathbf{W}\mathbf{W}^\top \mathbf{y} + \boldsymbol{\Sigma}_{\theta,0,\cdot} \boldsymbol{\Sigma}_\theta^{-1} \mathbf{P}\mathbf{P}^\top \mathbf{y}, \quad (6.15)$$

which is a linear mapping of the vector $(\mathbf{P}^\top \mathbf{H}\boldsymbol{\beta}, \mathbf{W}^\top \mathbf{y}, \mathbf{P}^\top \mathbf{y})^\top$. Now, Equation 6.12 tells us that conditional to $\mathbf{P}^\top \mathbf{y}$, $(\mathbf{P}^\top \mathbf{H}\boldsymbol{\beta}, \mathbf{W}^\top \mathbf{y}, \mathbf{P}^\top \mathbf{y})^\top$ is a (degenerate) Gaussian vector, so Gaussian theory implies that conditional to $\mathbf{P}^\top \mathbf{y}$, $(\mathbf{y}_0, \mathbf{P}^\top \mathbf{H}\boldsymbol{\beta}, \mathbf{W}^\top \mathbf{y}, \mathbf{P}^\top \mathbf{y})^\top$ is a Gaussian vector and therefore $(\mathbf{y}_0, \mathbf{P}^\top \mathbf{H}\boldsymbol{\beta} - \mathbf{P}^\top \mathbf{y}, -\mathbf{W}^\top \mathbf{y})^\top$ is one as well. So all that remains to be shown is that its mean and covariance are those given by Proposition 6.1.

To do this, we compute $\mathbf{E}_\theta^{\mathbf{y},\boldsymbol{\beta}}$ and $\sigma^2 \mathbf{S}_\theta^{\mathbf{y},\boldsymbol{\beta}}$, the conditional mean and variance of \mathbf{y}_0 given $\mathbf{W}^\top \mathbf{y}$, $\mathbf{P}^\top \mathbf{y}$ and $\mathbf{P}^\top \mathbf{H}\boldsymbol{\beta}$ and check that they fit the parameters of (6.13).

$$\begin{aligned} \mathbf{E}_\theta^{\mathbf{y},\boldsymbol{\beta}} &= \mathbf{E}_0 + \mathbf{S}_{\theta,0,\cdot} \begin{pmatrix} \mathbf{P}^\top \\ \mathbf{W}^\top \end{pmatrix} \boldsymbol{\Sigma}_\theta^{-1} \begin{pmatrix} \mathbf{P} & \mathbf{W} \end{pmatrix} \begin{pmatrix} \mathbf{P}^\top \mathbf{H}\boldsymbol{\beta} - \mathbf{P}^\top \mathbf{y} \\ -\mathbf{W}^\top \mathbf{y} \end{pmatrix} \\ &= \mathbf{H}_{0,0} \boldsymbol{\beta} + \boldsymbol{\Sigma}_{\theta,0,\cdot} \boldsymbol{\Sigma}_\theta^{-1} (\mathbf{y} - \mathbf{H}\boldsymbol{\beta}); \end{aligned} \quad (6.16)$$

$$\begin{aligned} \mathbf{S}_\theta^{\mathbf{y},\boldsymbol{\beta}} &= \mathbf{S}_{\theta,0,0} - \mathbf{S}_{\theta,0,\cdot} \begin{pmatrix} \mathbf{P}^\top \\ \mathbf{W}^\top \end{pmatrix} \boldsymbol{\Sigma}_\theta^{-1} \begin{pmatrix} \mathbf{P} & \mathbf{W} \end{pmatrix} \mathbf{S}_{\theta,\cdot,0} \\ &= \boldsymbol{\Sigma}_{\theta,0,0} - \boldsymbol{\Sigma}_{\theta,0,\cdot} \boldsymbol{\Sigma}_\theta^{-1} \boldsymbol{\Sigma}_{\theta,\cdot,0}. \end{aligned} \quad (6.17)$$

□

From this point onwards, distributions are no longer implicitly conditional to σ^2 and $\boldsymbol{\theta}$.

Corollary 6.2. *Assume that $\mathbf{P}\mathbf{P}^\top + \mathbf{W}\mathbf{W}^\top = \mathbf{I}_n$. The predictive distribution when $\boldsymbol{\beta}$ is unknown – i.e. the distribution of \mathbf{y}_0 conditional to \mathbf{y} , σ^2 and $\boldsymbol{\theta}$ – is Normal. With the notations of Proposition 6.1, it has mean vector $\mathbf{E}_0 - \mathbf{S}_{\theta,\cdot,0} \mathbf{W} (\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{y}$ and covariance matrix*

$$\sigma^2 \left\{ \mathbf{S}_{\theta,0,0} - \mathbf{S}_{\theta,\cdot,0} \mathbf{W} (\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{S}_{\theta,0,\cdot} \right\}.$$

Corollary 6.3. *Assume that $\mathbf{P}\mathbf{P}^\top + \mathbf{W}\mathbf{W}^\top = \mathbf{I}_n$. The predictive distribution when both β and σ^2 are unknown – i.e. the distribution of \mathbf{y}_0 conditional to \mathbf{y} and θ – is multivariate Student with $n - p$ degrees of freedom. With the notations of Proposition 6.1, it has location vector $\mathbf{E}_0 - \mathbf{S}_{\theta, \cdot, 0} \mathbf{W} \left(\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} \right)^{-1} \mathbf{W}^\top \mathbf{y}$ and scale matrix*

$$\frac{\mathbf{y}^\top \mathbf{W} \left(\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} \right)^{-1} \mathbf{W}^\top \mathbf{y}}{n - p} \left\{ \mathbf{S}_{\theta, 0, 0} - \mathbf{S}_{\theta, \cdot, 0} \mathbf{W} \left(\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} \right)^{-1} \mathbf{W}^\top \mathbf{S}_{\theta, 0, \cdot} \right\}.$$

6.3 Reference prior on a one-dimensional θ

In this section, θ is assumed to be a scalar parameter, which we emphasize by writing it θ . This is the easy case: Propositions 6.4 and 6.5 recall Proposition 3.2.

Proposition 6.4. *The reference prior on θ is $\pi(\theta) \propto$*

$$\sqrt{\text{Tr} \left[\left\{ \mathbf{W}^\top \frac{\partial}{\partial \theta} (\boldsymbol{\Sigma}_\theta) \mathbf{W} \left(\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} \right)^{-1} \right\}^2 \right] - \frac{1}{n - p} \text{Tr} \left\{ \mathbf{W}^\top \frac{\partial}{\partial \theta} (\boldsymbol{\Sigma}_\theta) \mathbf{W} \left(\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} \right)^{-1} \right\}^2}. \quad (6.18)$$

This result is in keeping with the previous work of [Berger et al., 2001]:

Proposition 6.5. *The reference prior on θ can also be written as:*

$$\pi(\theta) \propto \sqrt{\text{Tr} \left[\left\{ \frac{\partial}{\partial \theta} (\boldsymbol{\Sigma}_\theta) \boldsymbol{\Sigma}_\theta^{-1} \mathbf{Q}_\theta \right\}^2 \right] - \frac{1}{n - p} \left[\text{Tr} \left\{ \frac{\partial}{\partial \theta} (\boldsymbol{\Sigma}_\theta) \boldsymbol{\Sigma}_\theta^{-1} \mathbf{Q}_\theta \right\} \right]^2}, \quad (6.19)$$

where $\mathbf{Q}_\theta := \mathbf{I}_n - \mathbf{H} \left(\mathbf{H}^\top \boldsymbol{\Sigma}_\theta^{-1} \mathbf{H} \right)^{-1} \mathbf{H}^\top \boldsymbol{\Sigma}_\theta^{-1}$.

6.4 The Gibbs reference posterior on a multi-dimensional θ

Definition

In the case of a multidimensional θ , reference prior theory gives a choice between 1) considering θ as a single parameter or 2) defining an ordering on the scalar parameters $\theta_1, \dots, \theta_r$. Both possibilities are unsatisfactory, albeit in different ways. Concerning 1), Jeffreys' prior is unsuited to dealing with multidimensional parameters [Robert et al., 2009] and besides, the posterior may be improper. Concerning 2), further integration of the likelihood (6.10) would be analytically intractable, even if it were possible to define a non-arbitrary ordering of the coordinates of θ .

We propose a quasi-posterior distribution based on the reference posterior of models where only one coordinate of θ is unknown. For any integer $i \in \llbracket 1, r \rrbracket$, we collectively denote by $\boldsymbol{\theta}_{-i}$ all coordinates of θ except the i -th: $\boldsymbol{\theta}_{-i} = (\theta_j)_{j \in \llbracket 1, r \rrbracket \setminus \{i\}}$.

Consider now $\pi_i(\theta_i | \boldsymbol{\theta}_{-i})$, the reference prior distribution on θ_i conditional to $\boldsymbol{\theta}_{-i}$ and the associated reference posterior distribution $\pi_i(\theta_i | \mathbf{y}, \boldsymbol{\theta}_{-i}) \propto L^1(\mathbf{y} | \theta) \pi_i(\theta_i | \boldsymbol{\theta}_{-i})$.

The conditional reference prior $\pi_i(\theta_i | \boldsymbol{\theta}_{-i})$ is given by $\pi_i(\theta_i | \boldsymbol{\theta}_{-i}) \propto$

$$\sqrt{\text{Tr} \left[\left\{ \mathbf{W}^\top \partial_{\theta_i} \boldsymbol{\Sigma}_\theta \mathbf{W} \left(\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} \right)^{-1} \right\}^2 \right] - \frac{1}{n - p} \text{Tr} \left\{ \mathbf{W}^\top \partial_{\theta_i} \boldsymbol{\Sigma}_\theta \mathbf{W} \left(\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} \right)^{-1} \right\}^2}. \quad (6.20)$$

Now consider the sequence of conditional posterior distributions $(\pi_i(\theta_i|\mathbf{y}, \boldsymbol{\theta}_{-i}))_{i \in \llbracket 1, r \rrbracket}$. These conditional distributions are incompatible in the sense that there exists no joint probability distribution $\pi(\boldsymbol{\theta}|\mathbf{y})$ which agrees with all of them. We may however define the Gibbs reference posterior as a compromise between the conditionals in this sequence. In Chapter 4 we provided theoretical foundation for what such a compromise could be. In the end, we showed it to be the stationary probability distribution of a Markovian kernel $P_{\mathbf{y}} : (0, +\infty)^r \times \mathcal{B}((0, +\infty)^r)$, where $\mathcal{B}((0, +\infty)^r)$ denotes the Borel algebra on $((0, +\infty)^r)$. $P_{\mathbf{y}}$ is defined by the following expression, where $\boldsymbol{\theta}^{(0)} \in (0, 1)^r$ and δ_t denotes the shifted Dirac measure $\delta(\cdot - t)$:

$$P_{\mathbf{y}}(\boldsymbol{\theta}^{(0)}, d\boldsymbol{\theta}) = \frac{1}{r} \sum_{i=1}^r \pi_i(\theta_i|\mathbf{y}, \boldsymbol{\theta}_{-i}^{(0)}) d\theta_i \delta_{\boldsymbol{\theta}_{-i}^{(0)}}(d\boldsymbol{\theta}_{-i}). \quad (6.21)$$

The goal of this section is to provide sufficient conditions for the existence (and thus, propriety) of this stationary probability distribution $\pi_G(\boldsymbol{\theta}|\mathbf{y})$ and to show that the Markov Chain Monte-Carlo (MCMC) algorithm based on the Markovian kernel $P_{\mathbf{y}}$ converges to it, that is, $P_{\mathbf{y}}$ is uniformly ergodic. This means that denoting by $P_{\mathbf{y}}^n$ the Markov kernel produced by n successive applications of $P_{\mathbf{y}}$ and by $\|\cdot\|_{TV}$ the total variation norm,

$$\lim_{n \rightarrow \infty} \sup_{\boldsymbol{\theta}^{(0)} \in (0, +\infty)^r} \|P_{\mathbf{y}}^n(\boldsymbol{\theta}^{(0)}, \cdot) - \pi_G(\cdot|\mathbf{y})\|_{TV} = 0. \quad (6.22)$$

In the following results, when we write that “ $P_{\mathbf{y}}$ is uniformly ergodic”, we mean that Equation (6.22) holds.

Existence

The results in this subsection deal with the following setting:

- The spatial domain is the unit cube $(0, 1)^r$ ($r > 0$).
- The mean function space \mathcal{F}_p has dimension $p \geq 0$.
- The Universal Kriging model uses a Matérn anisotropic geometric or tensorized correlation kernel with smoothness parameter $\nu > 0$.
- Design sets contain $n > 0$ points, so we identify $(0, 1)^{rn}$ with the set of all design sets in the spatial domain $(0, 1)^r$. Let $Q(r, n)$ be the Lebesgue measure on $(0, 1)^{rn}$.

In the following, we change parametrization for the sake of convenience: define $\boldsymbol{\mu}$ such that $\forall i \in \llbracket 1, r \rrbracket$, $\mu_i = 1/\theta_i$. The conditionals are invariant to such a change, and therefore both the Markovian kernel $P_{\mathbf{y}}$ and, if it exists, its stationary probability remain the same. Abusing notations, the likelihood $L^1(\mathbf{y}|\boldsymbol{\theta})$ is denoted by $L(\mathbf{y}|\boldsymbol{\mu})$ when expressed in the $\boldsymbol{\mu}$ -parametrization. Define the functions f_i by $f_i(\mu_i|\boldsymbol{\mu}_{-i}) :=$

$$\sqrt{\text{Tr} \left[\left(\mathbf{W}^\top \frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_\mu \mathbf{W} \left(\mathbf{W}^\top \boldsymbol{\Sigma}_\mu \mathbf{W} \right)^{-1} \right)^2 \right] - \frac{1}{n-p} \text{Tr} \left[\mathbf{W}^\top \frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_\mu \mathbf{W} \left(\mathbf{W}^\top \boldsymbol{\Sigma}_\mu \mathbf{W} \right)^{-1} \right]^2}. \quad (6.23)$$

Then, following Equation (6.20), the conditional density π_i is in the $\boldsymbol{\mu}$ -parametrization given by:

$$\pi_i(\mu_i|\boldsymbol{\mu}_{-i}) \propto f_i(\mu_i|\boldsymbol{\mu}_{-i}). \quad (6.24)$$

We need to make some assumptions which are detailed below.

Assumption 3. *Any vector in the subspace of \mathbb{R}^n spanned by \mathbf{H} is either null or has strictly more than $2r$ non-null elements when expressed in the canonical base.*

Remark. It is not apparent, but the purpose of Assumption 3 is to control the behavior of the $f_i(\mu_i|\mu_{-i})$ ($i \in \llbracket 1, r \rrbracket$) when $\|\mu\| \rightarrow \infty$. See the proofs in Appendix 6.B for details.

This assumption is not very restrictive, as the two following results show.

Proposition 6.6. *In Ordinary Kriging – that is with $p = 1$ and \mathcal{F}_p being the space of constant functions – if $n > 2r$, Assumption 3 is automatically verified.*

Proof. In this setting, \mathbf{H} is a non-null constant $n \times 1$ matrix, so Assumption 3 is trivially verified. \square

Proposition 6.7. *Assume that the design set is such that any subset with cardinal $r+1$ forms a simplex. Then in Universal Kriging, if the mean function space \mathcal{F}_p is included within the vector space of polynomials of degree 0 and 1, and if $n > 3r$, Assumption 3 is automatically verified.*

Proof. Let \mathbf{y}^* belong to the subspace of \mathbb{R}^n spanned by \mathbf{H} . Assume that it has $2r$ or fewer non-null elements when expressed in the canonical base. Conversely, it has at least $n - 2r$ null elements. If $n > 3r$, then this means that there exists a function $f^* \in \mathcal{F}_p$ (the one represented by \mathbf{y}^*) which admits at least $r + 1$ zeros on the design set. However, given the premise of Proposition 6.7, these $r + 1$ points form a simplex, so they span an affine space of dimension r . As f^* is a polynomial with r unknowns of degree 0 or 1, this implies that $f^* = 0$. \square

Remark. $Q(r, n)$ -almost all design sets fit the premise of Proposition 6.7.

In some cases, Assumption 3 is sufficient for our purposes. Define $\mathbf{1}$ as the vector of \mathbb{R}^n with all components in the canonical basis equal to 1.

Proposition 6.8. *In the setting described above, if $0 < \nu < 1$ and $n > p + 1$, then for $Q(r, n)$ -almost all design sets, if $\mathbf{1}$ does not belong to the vector space spanned by \mathbf{H} , then Assumption 3 implies that there exists a hyperplane \mathcal{H} of \mathbb{R}^n such that $\forall \mathbf{y} \in \mathbb{R}^n \setminus \mathcal{H}$, $P_{\mathbf{y}}$ is uniformly ergodic.*

The proof of this Proposition can be found in Appendix 6.B.

Naturally, the above result is somewhat unsatisfactory since most users will want to include non-null constant functions in \mathcal{F}_p .

Consider now the following assumption.

Assumption 4. *There exists $\epsilon_{\mathbf{y}} > 0$ such that $L(\mathbf{y}|\mu) = O(\|\mu\|^{\epsilon_{\mathbf{y}}})$ when $\|\mu\| \rightarrow 0$.*

Remark. Assumption 4 essentially means that the model should find perfect correlation unlikely.

The following theorem, which is proved in Appendix 6.B, is our essential tool for dealing with the case where non-null constant functions are included in \mathcal{F}_p .

Theorem 6.9. *In the setting described above, if $\nu > 1$ and $n > p + 2r + 2$, then for $Q(r, n)$ -almost all design sets, Assumptions 3 and 4 imply that $P_{\mathbf{y}}$ is uniformly ergodic.*

The next two results, which are proved in Appendix 6.B, concern particular settings where Assumptions 3 and 4 are both verified and therefore Theorem 6.9 yields the uniform ergodicity of $P_{\mathbf{y}}$.

Proposition 6.10. *Consider the particular case of the above described setting where $p = 1$ and \mathcal{F}_p is the space of all constant functions (Ordinary Kriging), and assume that one of the following conditions is satisfied:*

1. $1 < \nu < 2$ and $n > 2r + 3$;
2. $2 < \nu < 3$ and $n > (r + 1)(r/2 + 2)$.

Then, for $Q(r, n)$ -almost all design sets, there exists a hyperplane \mathcal{H} of \mathbb{R}^n such that $\forall \mathbf{y} \in \mathbb{R}^n \setminus \mathcal{H}$, $P_{\mathbf{y}}$ is uniformly ergodic.

Proposition 6.11. *Consider the particular case of the above described setting where \mathcal{F}_p is included within the space of all polynomials of degree 0 and 1 (so $p \leq r + 1$) and assume that the following condition is satisfied:*

- $2 < \nu < 3$ and $n > r(r + 1)/2 + 2r + 3$.

Then, for $Q(r, n)$ -almost all design sets, there exists a hyperplane \mathcal{H} of \mathbb{R}^n such that $\forall \mathbf{y} \in \mathbb{R}^n \setminus \mathcal{H}$, $P_{\mathbf{y}}$ is uniformly ergodic.

Remark. In Propositions 6.8, 6.10 and 6.11, the condition that the observation \mathbf{y} should not belong to a given negligible (for the Lebesgue measure) subset of \mathbb{R}^n is fairly natural: for the Kriging model to be adequate, \mathbf{y} must not look like a realization of a degenerate Gaussian vector. Theorem 6.9 does not really dispense with it, as it is implied by Assumption 4.

To sum up the results of this section, to ensure that the Gibbs reference posterior exists and can be accessed through Gibbs sampling, one should check that one of the following assertions is true:

- $\nu > 1$, $n > r + p + 2$ and both Assumptions 3 and 4 are verified;
- \mathcal{F}_p contains only constant functions and $1 < \nu < 2$ and $n > r + 3$;
- \mathcal{F}_p contains only polynomials of degree 0 and 1, $2 < \nu < 3$ and $n > r(r + 1)/2 + 2r + 3$;
- $0 < \nu < 1$, $n > p + 1$, no non-null constant function belongs to \mathcal{F}_p and Assumption 3 is verified.

6.5 Comparison of the predictive performance of the full-Bayesian approach versus MLE and MAP plug-in approaches

In this section, we evaluate the predictive performance resulting from the Gibbs reference posterior distribution $\pi_G(\boldsymbol{\theta}|\mathbf{y})$ in the context of a well-specified model, and then when emulating some deterministic real functions. We contrast the full-Bayesian approach, in which the Full Gibbs reference Posterior Distribution (FPD) is used, with two plug-in approaches: one where the Maximum Likelihood Estimator (MLE) and the other where the Maximum A Posteriori (MAP) estimator is assumed to be the true value of $\boldsymbol{\theta}$. All approaches make use of the reference posterior $\pi(\boldsymbol{\beta}, \sigma^2|\mathbf{y}, \boldsymbol{\theta})$.

We use the following terminology. We call Simple Kriging the Kriging model where the mean function is assumed to be known, whether this assumption is correct or known. We call Ordinary Kriging any Universal Kriging model where the mean function space is the space

of constant functions. We call Affine Kriging any Universal Kriging model where the mean function space is the space of affine functions.

Well-specified model

We first consider well-specified models, specifically Kriging models with unknown parameters $(\beta, \sigma^2, \theta)$ emulating actual Gaussian processes with variance $\sigma^2 = 1$ and Matérn anisotropic geometric autocorrelation kernel with smoothness $\nu = 5/2$. Moreover, the true mean function of the Gaussian process belongs to the assumed mean function space \mathcal{F}_p .

The spatial domain is the unit cube $(0, 1)^r$ and the considered design sets all contain n points independently chosen according to the Lebesgue measure on the domain $(0, 1)^r$.

The following tables give the average coverage and average mean length of prediction intervals. To define these notions, we introduce the following notations:

- \mathbf{Y} is the Gaussian process, and $\mathbf{Y}(\mathbf{x})$ is the vector of the values taken by said process at the points in the design set \mathbf{x} ;
- T is a random variable which follows the Uniform distribution on the unit cube $(0, 1)^r$. It represents the “test” point;
- \mathbf{X} is the random design set following the Uniform distribution on $((0, 1)^r)^n$.
- \mathbf{Y} , T and \mathbf{X} are mutually independent;
- f is a function defined on $((0, 1)^r)^n \times \mathbb{R}^n \times (0, 1)^r$ which associates to $(\mathbf{x}, \mathbf{y}, t)$ the prediction interval at t of the Gaussian process, based on the knowledge of its value \mathbf{y} on the design set \mathbf{x} .

Definition 6.12. *The average coverage is the probability (with respect to the distributions of \mathbf{X} , \mathbf{Y} and T) that $\mathbf{Y}(T) \in f(\mathbf{X}, \mathbf{Y}(\mathbf{X}), T)$.*

Definition 6.13. *The average mean length is the expectation (with respect to the distributions of \mathbf{X} , \mathbf{Y} and T) of the length of $f(\mathbf{X}, \mathbf{Y}(\mathbf{X}), T)$.*

The average coverage $\mathbb{P}[\mathbf{Y}(T) \in f(\mathbf{X}, \mathbf{Y}(\mathbf{X}), T)]$ is numerically computed as

$$\mathbb{P}[\mathbf{Y}(T) \in f(\mathbf{X}, \mathbf{Y}(\mathbf{X}), T)] = \mathbb{E}[\mathbb{P}[\mathbf{Y}(T) \in f(\mathbf{X}, \mathbf{Y}(\mathbf{X}), T) | \mathbf{X}, \mathbf{Y}(\mathbf{X})]]$$

over 500 random design sets and for each design set 1000 random test points. The average mean length is computed in a similar fashion.

In this subsection and the following one, we take $n = 30$ and $r = 3$.

In the first set of simulations, we use a well-specified Ordinary Kriging model, with the unknown mean 5. As $r = 3$ and $n = 30$, Proposition 6.10 is applicable.

The results given in Table 6.1 show that using the full posterior distribution (FPD) to derive the predictive distribution is the best possible choice from a frequentist point of view as the nominal value is nearly matched by the average coverage. Predictive Intervals derived from the MAP estimator do not perform as well, and Predictive Intervals derived from the MLE perform even worse.

The results given in Table 6.2 show that Predictive Intervals arising from the full Gibbs reference posterior distribution (FPD) are on average somewhat larger than those resulting

Average Coverage				
Corr. lengths	True	MLE	MAP	FPD
0.4 – 0.8 – 0.2	0.95	0.88	0.91	0.95
0.5 – 0.5 – 0.5	0.95	0.88	0.90	0.94
0.7 – 1.3 – 0.4	0.95	0.90	0.92	0.95
0.8 – 0.3 – 0.6	0.95	0.89	0.91	0.94
0.8 – 1.0 – 0.9	0.95	0.90	0.92	0.94

Table 6.1 – For a Gaussian Process with constant mean function equal to 5, variance parameter 1 and smoothness parameter 5/2, average coverage of 95% Prediction Intervals produced by an Ordinary Kriging model. “True” stands for the Simple Kriging prediction based on the knowledge of the true mean parameter, variance parameter and vector of correlation lengths.

from knowledge of the true parameters, while intervals arising from both types of parameter estimation (MLE and MAP) are too short.

Average Mean Length				
Corr. lengths	True	MLE	MAP	FPD
0.4 – 0.8 – 0.2	2.23	2.06	2.14	2.58
0.5 – 0.5 – 0.5	1.69	1.55	1.59	1.83
0.7 – 1.3 – 0.4	1.09	1.02	1.07	1.20
0.8 – 0.3 – 0.6	1.63	1.51	1.57	1.81
0.8 – 1.0 – 0.9	0.71	0.66	0.69	0.76

Table 6.2 – For a Gaussian Process with constant mean function equal to 5, variance parameter 1 and smoothness parameter 5/2, average mean length of 95% Prediction Intervals produced by an Ordinary Kriging model. “True” stands for the Simple Kriging prediction based on the knowledge of the true mean parameter, variance parameter and vector of correlation lengths.

Consider now Universal Kriging models where the true mean function is the polynomial $(x_1, x_2, x_3) \mapsto 5 + 4x_1 + 3x_2 + 2x_3$, and the model (correctly) assumes that it belongs to the 4-dimensional space ($p = 4$) spanned by the functions mapping (x_1, x_2, x_3) to 1, x_1 , x_2 and x_3 respectively. For such Affine Kriging models, Proposition 6.11 is applicable.

As shown in Table 6.3, Predictive Intervals resulting from both plug-in approaches (MLE, MAP) and from the full posterior distribution perform a little worse than in the Ordinary Kriging setting, but their relative performances stay the same.

Table 6.4 shows that the average mean lengths of Predictive Intervals are not very different in Affine Kriging than in Ordinary Kriging when it comes to the FPD. However, they are

Average Coverage				
Corr. lengths	True	MLE	MAP	FPD
0.4 – 0.8 – 0.2	0.95	0.87	0.90	0.94
0.5 – 0.5 – 0.5	0.95	0.87	0.89	0.92
0.7 – 1.3 – 0.4	0.95	0.89	0.92	0.94
0.8 – 0.3 – 0.6	0.95	0.87	0.90	0.93
0.8 – 1.0 – 0.9	0.95	0.89	0.92	0.93

Table 6.3 – For a Gaussian Process with mean function $(x_1, x_2, x_3) \mapsto 5 + 4x_1 + 3x_2 + 2x_3$, variance parameter 1 and smoothness parameter 5/2, average coverage of 95% Prediction Intervals produced by an Affine Kriging model. “True” stands for the Simple Kriging prediction based on the knowledge of the true mean function, variance parameter and vector of correlation lengths.

larger in Affine Kriging than in Ordinary Kriging when it comes to the MLE and the MAP. Interestingly, Predictive Intervals resulting from the MAP have about the same size as Predictive Intervals derived when all parameters are known. Those derived using the MLE are shorter, and those derived from the FPD are larger.

Average Mean Length				
Corr. lengths	True	MLE	MAP	FPD
0.4 – 0.8 – 0.2	2.23	2.14	2.23	2.59
0.5 – 0.5 – 0.5	1.69	1.57	1.66	1.83
0.7 – 1.3 – 0.4	1.09	1.04	1.10	1.20
0.8 – 0.3 – 0.6	1.63	1.54	1.61	1.80
0.8 – 1.0 – 0.9	0.71	0.67	0.71	0.75

Table 6.4 – For a Gaussian Process with mean function $(x_1, x_2, x_3) \mapsto 5 + 4x_1 + 3x_2 + 2x_3$, variance parameter 1 and smoothness parameter $5/2$, average mean length of 95% Prediction Intervals produced by an Affine Kriging model. “True” stands for the Simple Kriging prediction based on the knowledge of the true mean function, variance parameter and vector of correlation lengths.

For reference, we give the tables obtained in the Simple Kriging case, that is the case where the Gaussian Process is known to have null mean function. Table 6.5 gives the average coverages and Table 6.6 the average mean lengths.

Average Coverage				
Corr. lengths	True	MLE	MAP	FPD
0.4 – 0.8 – 0.2	0.95	0.88	0.91	0.95
0.5 – 0.5 – 0.5	0.95	0.89	0.90	0.94
0.7 – 1.3 – 0.4	0.95	0.90	0.92	0.95
0.8 – 0.3 – 0.6	0.95	0.89	0.91	0.95
0.8 – 1.0 – 0.9	0.95	0.90	0.92	0.94

Table 6.5 – For a Gaussian Process with null mean function, variance parameter 1 and smoothness parameter $5/2$, average coverage of 95% Prediction Intervals produced by a Simple Kriging model. “True” stands for the prediction based on the knowledge of the true variance parameter and the true vector of correlation lengths.

Average Mean Length				
Corr. lengths	True	MLE	MAP	FPD
0.4 – 0.8 – 0.2	2.23	2.05	2.13	2.59
0.5 – 0.5 – 0.5	1.69	1.55	1.58	1.84
0.7 – 1.3 – 0.4	1.09	1.02	1.07	1.21
0.8 – 0.3 – 0.6	1.63	1.51	1.56	1.82
0.8 – 1.0 – 0.9	0.71	0.66	0.69	0.76

Table 6.6 – For a Gaussian Process with null mean function, variance parameter 1 and smoothness parameter $5/2$, average mean length of 95% Prediction Intervals produced by a Simple Kriging model. “True” stands for the prediction based on the knowledge of the true variance parameter and the true vector of correlation lengths.

The performance of Ordinary Kriging when the mean function is constant is nearly the same as that of Simple Kriging when the mean function is known.

The performance of Affine Kriging when the mean function is affine, however, is noticeably poorer than the performance of Simple Kriging when the mean function is known: its average coverage is lower. This is not too surprising, since the prediction problem is more difficult.

Average Coverage				
Corr. lengths	True	MLE	MAP	FPD
0.4 – 0.8 – 0.2	0.95	0.77	0.81	0.88
0.5 – 0.5 – 0.5	0.95	0.80	0.82	0.89
0.7 – 1.3 – 0.4	0.95	0.82	0.86	0.91
0.8 – 0.3 – 0.6	0.95	0.79	0.83	0.89
0.8 – 1.0 – 0.9	0.95	0.82	0.86	0.91

Table 6.7 – For a Gaussian Process with mean function $(x_1, x_2, x_3) \mapsto 5 + 4x_1 + 3x_2 + 2x_3$, variance parameter 1 and smoothness parameter $5/2$, average coverage of 95% Prediction Intervals resulting from Simple Kriging (assuming the mean function is null for MLE/MAP/FPD and knowing it is $(x_1, x_2, x_3) \mapsto 5 + 4x_1 + 3x_2 + 2x_3$ for “True”). “True” stands for the prediction based on the knowledge of the true mean function, variance parameter and vector of correlation lengths.

Misspecified models

In this subsection, we deal with the performance of Kriging in cases where the Gaussian Process does not fit all assumptions.

First, we evaluate the performance of Universal Kriging in a context where the true mean function does not belong to the assumed mean function space \mathcal{F}_p . Precisely, we consider a Gaussian process with mean function $(x_1, x_2, x_3) \mapsto 5 + 4x_1 + 3x_2 + 2x_3$ and evaluate the performance of Simple Kriging (assuming the mean function to be null) with respect to that of Affine Kriging, which is the correct model in this situation.

Tables 6.7 and 6.8 show that Simple Kriging performs significantly worse than Affine Kriging when the mean function is $(x_1, x_2, x_3) \mapsto 5 + 4x_1 + 3x_2 + 2x_3$, both in terms of average coverage and average mean length of Predictive Intervals. Relative performances of MLE, MAP and FPD once again stay the same, though.

Average Mean Length				
Corr. lengths	True	MLE	MAP	FPD
0.4 – 0.8 – 0.2	2.23	2.23	2.36	2.78
0.5 – 0.5 – 0.5	1.69	1.61	1.66	1.92
0.7 – 1.3 – 0.4	1.09	1.03	1.13	1.28
0.8 – 0.3 – 0.6	1.63	1.54	1.63	1.87
0.8 – 1.0 – 0.9	0.71	0.64	0.68	0.77

Table 6.8 – For a Gaussian Process with mean function $(x_1, x_2, x_3) \mapsto 5 + 4x_1 + 3x_2 + 2x_3$, variance parameter 1 and smoothness parameter $5/2$, average mean length of 95% Prediction Intervals resulting from Simple Kriging (assuming the mean function is null for MLE/MAP/FPD and knowing it is $(x_1, x_2, x_3) \mapsto 5 + 4x_1 + 3x_2 + 2x_3$ for “True”). “True” stands for the prediction based on the knowledge of the true mean function, variance parameter and vector of correlation lengths.

This observation may lead us to investigate how Simple Kriging behaves with respect to Affine Kriging when the Gaussian Process is smoother than expected. Table 6.9 gives the average coverage and average mean length of Prediction Intervals resulting from the same procedure as before – that is, the correlation kernel is assumed to be Matérn with smoothness $5/2$ – but the Gaussian Process actually has a Squared Exponential correlation kernel (with correlation lengths 0.4, 0.8 and 0.2). These results can be compared with those from Table 6.10, which gives the results obtained when both the actual and the assumed correlation kernel are Matérn with smoothness $5/2$ (and the true correlation lengths are also 0.4, 0.8 and

0.2). It is apparent that performance is better when the actual kernel is Squared Exponential, both in terms of average coverage and average mean length. Recalling that this kernel can be seen as the limit of the Matérn kernel when the smoothness parameter goes to infinity, we conclude that a smoother process leads to an increase in performance for Simple, Ordinary and Affine Kriging. For Affine Kriging, the smoother process makes Prediction Intervals on average shorter, while the average coverage remains about the same. For Simple Kriging and to a lesser degree Ordinary Kriging, the smoother process makes Prediction Intervals on average shorter, while also increasing average coverage.

SQUARED EXPONENTIAL KERNEL	Average coverage			Average mean length		
Kriging model	MLE	MAP	FPD	MLE	MAP	FPD
Simple Kriging*	0.83	0.86	0.92	1.63	1.76	2.02
Ordinary Kriging	0.88	0.90	0.93	1.70	1.79	2.01
Affine Kriging	0.89	0.91	0.93	1.63	1.70	1.88

Table 6.9 – For a Gaussian Process with mean function $(x_1, x_2, x_3) \mapsto 5 + 4x_1 + 3x_2 + 2x_3$, variance parameter 1, and squared exponential correlation kernel with correlation lengths 0.4 - 0.8 - 0.2, average coverage and average mean length of 95% Prediction Intervals resulting from different types of Kriging (assuming the smoothness parameter to be 5/2). *With Simple Kriging, the mean function is (wrongly here) assumed to be constant with null value.

MATÉRN KERNEL $\nu = 5/2$	Average coverage			Average mean length		
Kriging model	MLE	MAP	FPD	MLE	MAP	FPD
Simple Kriging*	0.77	0.81	0.88	2.23	2.36	2.77
Ordinary Kriging	0.84	0.86	0.91	2.30	2.37	2.71
Affine Kriging	0.87	0.90	0.94	2.14	2.23	2.59

Table 6.10 – For a Gaussian Process with mean function $(x_1, x_2, x_3) \mapsto 5 + 4x_1 + 3x_2 + 2x_3$, variance parameter 1 and Matérn kernel with correlation lengths 0.4 - 0.8 - 0.2 and smoothness parameter 5/2, average coverage and average mean length of 95% Prediction Intervals resulting from different types of Kriging. *With Simple Kriging, the mean function is (wrongly here) assumed to be constant with null value.

All else being equal, smoother processes result in a better quality of prediction for Simple, Ordinary and Affine Kriging, because the observed values of the process yield more information about the value of the process in the neighborhoods of the observation points. This even makes up to some degree for the misspecification of the mean function, so the improvement is greater in the case of Simple Kriging.

Emulating deterministic functions

In this subsection, we test the ability of the model to predict deterministic functions, namely the 7-dimensional Ackley and Rastrigin functions. The Ackley and the Rastrigin functions have the following expressions:

$$A(\mathbf{x}) = 20 + \exp(1) - 20 \exp \left(-0.2 \sqrt{\frac{1}{7} \sum_{i=1}^7 x_i^2} \right) - \exp \left(\frac{1}{7} \sum_{i=1}^7 \cos(2\pi x_i) \right); \quad (6.25)$$

$$R(\mathbf{x}) = 70 + \sum_{i=1}^7 (x_i^2 - 10 \cos(2\pi x_i)). \quad (6.26)$$

Naturally, the notions of average coverage and average mean length for Prediction intervals make no sense in this setting, since we can no longer average our results over the distribution of a Gaussian process. Denoting d the deterministic function, and using previous notations, we may define:

Definition 6.14. *The coverage is the probability (with respect to the distribution of \mathbf{X} and T) that $d(T) \in f(\mathbf{X}, d(\mathbf{X}), T)$.*

Definition 6.15. *The mean length is the expectation (with respect to the distribution of \mathbf{X} and T) of the length of $f(\mathbf{X}, d(\mathbf{X}), T)$.*

The coverage $\mathbb{P}[\mathbf{Y}(T) \in f(\mathbf{X}, d(\mathbf{X}), T)]$ is numerically computed as

$$\mathbb{P}[\mathbf{Y}(T) \in f(\mathbf{X}, \mathbf{Y}(\mathbf{X}), T)] = \mathbb{E}[\mathbb{P}[\mathbf{Y}(T) \in f(\mathbf{X}, d(\mathbf{X}), T) | \mathbf{X}]]$$

over 500 design sets and for each design set 1000 test points. The mean length is computed in a similar fashion.

When emulating the Ackley or the Rastrigin function, we take $r = 7$ and $n = 100$.

We must stress that there is no reason that the coverage of 95% Prediction Intervals, whether produced by MLE or MAP plug-in methods or by the full Gibbs reference posterior distribution should be 95%, but depending on whether or not Kriging can be considered a good surrogate model for the Ackley or Rastrigin function, the coverage of 95% Prediction Intervals may be more or less close to the 95% target figure.

First, we consider an Ordinary Kriging model with anisotropic geometric Matérn kernel of smoothness $\nu = 5/2$. Because $r = 7$ and $n = 100$, Proposition 6.10 is applicable.

When emulating the Ackley function (cf. Table 6.11), regardless of the Kriging method used, the full posterior distribution significantly improves the average coverage of Prediction Intervals when compared to the MLE or the MAP, with a comparatively small trade-off regarding the mean length of these intervals. This result is consistent with results obtained with actual realizations of Gaussian processes.

When we emulate the Rastrigin function (cf. Table 6.12), coverages come closer to the average coverages given in Tables 6.1, 6.3 and 5.2. But the more significant fact of the improvement of the coverage by the full posterior distribution is as true here as in the Ackley case. We may simply infer from this that the Rastrigin function can more plausibly be seen as a realization of a Gaussian Process than the Ackley function.

Let us now compare the performance of different Kriging models: Simple (mean function assumed null), Ordinary and Affine. When emulating the Ackley function, Ordinary and Affine Kriging models yield slightly higher Prediction Interval coverages than Simple Kriging,

EMULATED FUNCTION: ACKLEY	Coverage			Mean length		
Kriging model	MLE	MAP	FPD	MLE	MAP	FPD
Simple Kriging	0.84	0.87	0.90	0.35	0.36	0.39
Ordinary Kriging	0.87	0.88	0.91	0.37	0.38	0.41
Affine Kriging	0.87	0.90	0.91	0.37	0.39	0.41

Table 6.11 – Coverage and mean length of 95% Prediction Intervals when emulating the 7-dimensional Ackley function (Matérn anisotropic geometric correlation kernel with smoothness $\nu = 5/2$).

EMULATED FUNCTION: RASTRIGIN	Coverage			Mean length		
Kriging model	MLE	MAP	FPD	MLE	MAP	FPD
Simple Kriging	0.94	0.94	0.96	28.3	28.3	30.2
Ordinary Kriging	0.91	0.92	0.94	26.2	26.7	28.3
Affine Kriging	0.90	0.91	0.92	25.9	26.5	27.2

Table 6.12 – Coverage and mean length of 95% Prediction Intervals when emulating the 7-dimensional Rastrigin function (Matérn anisotropic geometric correlation kernel with smoothness $\nu = 5/2$).

at the cost of slightly higher mean lengths. When emulating the Rastrigin function, we actually observe the reverse phenomenon.

From this study, we cannot conclusively ascertain whether Universal Kriging, at least in the form of Ordinary or Affine Kriging, yields better results than Simple Kriging. All that can be said is that these Kriging methods are more or less conservative, but even this depends on the emulated function.

In the following example (cf. Table 6.13), we add the linear function $(x_1, x_2, x_3, x_4, x_5, x_6, x_7) \mapsto 100 \sum_{i=1}^7 x_i$ to the 7-dimensional Rastrigin function. We may expect this modification of the Rastrigin function to be more accurately emulated by Affine Kriging than by Simple Kriging.

RASTRIGIN + $100 \sum_{i=1}^7 x_i$	Coverage			Mean length		
Kriging model	MLE	MAP	FPD	MLE	MAP	FPD
Simple Kriging	0.88	0.92	0.94	25.9	29.3	31.1
Ordinary Kriging	0.87	0.91	0.93	25.7	28.5	30.4
Affine Kriging	0.90	0.91	0.92	26.0	26.6	27.3

Table 6.13 – Coverage and mean length of 95% Prediction Intervals when emulating the 7-dimensional Rastrigin function augmented by a linear function (Matérn anisotropic geometric correlation kernel with smoothness $\nu = 5/2$).

The addition of the linear function causes a decrease in performance for Prediction Intervals of both Simple and Ordinary Kriging, in the sense that coverage decreases while mean length increases for MAP and FPD. And the coverage of MLE sinks so much – from 94% to 88% for Simple Kriging and from 91% to 87% for Ordinary Kriging – that its performance may also be said to decrease, even though its mean length is slightly lower.

The performance of Affine Kriging is unchanged, however, whether one considers the MLE or MAP plug-in methods or the method using the full posterior distribution. This suggests

that with a stronger linear component, Affine Kriging would be clearly preferable to Simple or Ordinary Kriging.

To test this, we emulate the 7-dimensional Rastrigin function, to which we add a stronger linear term: $(x_1, x_2, x_3, x_4, x_5, x_6, x_7) \mapsto 120 \sum_{i=1}^7 x_i$.

RASTRIGIN + $120 \sum_{i=1}^7 x_i$	Coverage			Mean length		
	MLE	MAP	FPD	MLE	MAP	FPD
Simple Kriging	0.90	0.94	0.96	27.6	33.0	37.6
Ordinary Kriging	0.88	0.92	0.94	26.9	30.5	32.3
Affine Kriging	0.90	0.92	0.92	26.0	26.9	27.3

Table 6.14 – Coverage and mean length of 95% Prediction Intervals when emulating the 7-dimensional Rastrigin function augmented by a linear function (Matérn anisotropic geometric correlation kernel with smoothness $\nu = 5/2$).

For Simple Kriging, Prediction Intervals coverage and mean length are both higher when $120 \sum_{i=1}^7 x_i$ is added to the Rastrigin function (Table 6.14) rather than $100 \sum_{i=1}^7 x_i$ (Table 6.13). This is also true, though to a lesser extent, of Ordinary Kriging. The performance of Affine Kriging, on the other hand, still remains the same because it can account for any linear term by seeing it as part of the mean function. Simple Kriging (assuming the mean function to be null) and Ordinary Kriging do not have this luxury and must assume a greater variance for the Gaussian process, which results in more conservative Predictive Intervals.

Gathering the results obtained above, we conclude that Universal Kriging only significantly improves performance if the trend belongs to the assumed mean function space \mathcal{F}_p and if it stands out. In other words, the signal/noise ratio must be high, where the signal is here the “true” mean function and the noise is the stationary Gaussian Process added to it. When no trend of the expected form can be discerned, like when emulating the Ackley or Rastrigin function through Affine Kriging, then there is no significant benefit to using Universal instead of Simple Kriging. When the ratio is high, as in the case of the Rastrigin function with the addition of the greater linear term $120 \sum_{i=1}^7 x_i$, Universal Kriging (if the mean function space \mathcal{F}_p is adequately defined) improves upon Simple Kriging, which becomes overly conservative. Further, when the emulated function is particularly smooth, Simple Kriging becomes capable of capturing the trend to some extent even if the mean function is misspecified, thanks to the mechanics of Gaussian conditioning.

6.6 Conclusion

In this chapter, we provided an Objective Bayesian solution to the problem of taking into account parameter uncertainty when performing prediction based on a Universal Kriging model with anisotropic Matérn autocorrelation kernel. The reference posterior on the location parameter β and the variance parameter σ^2 is coupled with the Gibbs reference posterior on the vector of correlation lengths θ . By using the Gibbs reference posterior, which is the optimal compromise between the conditional reference posteriors on one correlation length θ_i based on the knowledge of all other correlation length θ_j ($j \neq i$), we bypass the problem of determining an ordering on the correlation lengths. Moreover, this solution allows for

Gibbs sampling of the posterior distribution, which makes full-Bayesian inference or prediction tractable.

We proved that the Gibbs reference posterior exists and is proper in several Universal Kriging settings, depending on the number of available observation points and on the smoothness parameter of the Matérn kernel.

Numerical simulations show that Prediction Intervals produced by the full-Bayesian procedure based on the Gibbs reference posterior have better coverage than those produced by the Maximum Likelihood Estimator or even the Maximum A Posteriori estimator, and that their mean length is only moderately greater.

In addition, these simulations showed that when emulating deterministic functions, there is no obvious advantage to using Universal Kriging over Simple Kriging, unless the trend strongly stands out and belongs to the assumed mean function space.

From a theoretical standpoint, the Universal Kriging setting poses specific problems when compared to the Simple Kriging setting. As was shown (to our knowledge for the first time) by Berger et al. [2001], the behavior of the integrated likelihood changes significantly depending on whether functions that take a non-null constant value on the design set are included in the mean function space \mathcal{F}_p . The integrated likelihood often fails to vanish in the neighborhood of perfect correlation in Ordinary Kriging models and *a fortiori* in more complex Universal Kriging models where the constant term of the mean function is unknown. Berger et al. [2001] show in the isotropic framework that the reference prior adapts to this situation by being proper (at least for sufficiently rough correlation kernels – we provided a proof for smoother kernels in Chapter 3). In the anisotropic framework however, we were not able to prove the existence of the Gibbs reference posterior in such a situation, which is why we require Assumption 4. Although it is possible that closer analysis may allow us to relax this requirement, we find it more likely that Assumption 4 is the price we pay for defining the Gibbs reference posterior as a compromise between incompatible conditional reference posterior distributions. Indeed, each conditional maximizes the expected information of the model when all but one correlation length are fixed at finite values, i.e. in a context where *perfect correlation is impossible*, whatever may be the value of the unfixed correlation length. Therefore, it is conceivable that in the absence of penalization by the integrated likelihood of the kind given by Assumption 4, the conditionals may place too much weight on high values of the unfixed correlation length for the Gibbs reference posterior to be well defined.

Taking this restriction into account in Theorem 6.9, we proved that the Gibbs reference posterior exists and is the limit of a uniformly converging Markov Chain Monte-Carlo (MCMC) algorithm for commonly used Matérn anisotropic geometric and tensorized correlation kernels when the design set has enough points (cf. Propositions 6.10 and 6.11). More generally, we would conjecture that for any noninteger smoothness $\nu \in (1, +\infty)$, and if the mean function space \mathcal{F}_p does not contain polynomials of degree higher than $[\nu] - 1$, there exists some lower bound on the cardinal of the design set over which the Gibbs reference posterior exists and the MCMC algorithm uniformly converges to it. However, this lower bound may be too high for practical purposes.

Future work may involve gaining a better understanding of the significance of the Gibbs reference posterior as a compromise between the incompatible reference conditionals on correlation

lengths. This method was primarily intended as a practical means of solving the problem of giving an objective posterior distribution on correlation lengths in the case of anisotropic correlation kernels, where the reference posterior is intractable and may not be proper. But its theoretical properties beyond its propriety, its invariance under reparametrizations of the type $f((\theta_1, \dots, \theta_r)^\top) = (f_1(\theta_1), \dots, f_r(\theta_r))^\top$ and its apparent good frequentist performances remain unknown.

Appendix 6.A Matérn kernels

We use the following convention for the Fourier transform: the Fourier transform \widehat{g} of a smooth function $g : \mathbb{R}^r \rightarrow \mathbb{R}$ verifies $g(\mathbf{x}) = \int_{\mathbb{R}^r} \widehat{g}(\boldsymbol{\omega}) e^{i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} d\boldsymbol{\omega}$ and $\widehat{g}(\boldsymbol{\omega}) = (2\pi)^{-r} \int_{\mathbb{R}^r} g(\mathbf{x}) e^{-i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} d\mathbf{x}$. Let us set up a few notations.

- (a) \mathcal{K}_ν is the modified Bessel function of second kind with parameter ν ;
 (b) $K_{r,\nu}$ is the r -dimensional Matérn isotropic covariance kernel with variance 1, correlation length 1 and smoothness $\nu \in (0, +\infty)$ and $\widehat{K}_{r,\nu}$ is its Fourier transform:

$$(i) \quad \forall \mathbf{x} \in \mathbb{R}^r, \quad K_{r,\nu}(\mathbf{x}) = \frac{1}{\Gamma(\nu)2^{\nu-1}} (2\sqrt{\nu}\|\mathbf{x}\|)^\nu \mathcal{K}_\nu(2\sqrt{\nu}\|\mathbf{x}\|) ; \quad (6.27)$$

$$(ii) \quad \forall \boldsymbol{\omega} \in \mathbb{R}^r, \quad \widehat{K}_{r,\nu}(\boldsymbol{\omega}) = \frac{M_r(\nu)}{(\|\boldsymbol{\omega}\|^2 + 4\nu)^{\nu + \frac{r}{2}}} \text{ with } M_r(\nu) = \frac{\Gamma(\nu + \frac{r}{2})(2\sqrt{\nu})^{2\nu}}{\pi^{\frac{r}{2}}\Gamma(\nu)}. \quad (6.28)$$

- (c) $K_{r,\nu}^{tens}$ is the r -dimensional Matérn tensorized covariance kernel with variance 1, correlation length 1 and smoothness $\nu \in \mathbb{R}_+$ and $\widehat{K}_{r,\nu}^{tens}$ is its Fourier transform:

$$(i) \quad \forall \mathbf{x} \in \mathbb{R}^r, \quad K_{r,\nu}^{tens}(\mathbf{x}) = \prod_{j=1}^r K_{1,\nu}(\mathbf{x}_j) ; \quad (6.29)$$

$$(ii) \quad \forall \boldsymbol{\omega} \in \mathbb{R}^r, \quad \widehat{K}_{r,\nu}^{tens}(\boldsymbol{\omega}) = \prod_{j=1}^r \widehat{K}_{1,\nu}(\boldsymbol{\omega}_j). \quad (6.30)$$

- (d) let us adopt the following convention: if $\mathbf{t} \in \mathbb{R}^r$, $\frac{\mathbf{t}}{\boldsymbol{\theta}} = \left(\frac{t_1}{\theta_1}, \dots, \frac{t_r}{\theta_r}\right)$ and $\mathbf{t}\boldsymbol{\mu} = (t_1\mu_1, \dots, t_r\mu_r)$.

We define the Matérn geometric anisotropic covariance kernel with variance parameter σ^2 , correlation lengths $\boldsymbol{\theta}$ (resp. inverse correlation lengths $\boldsymbol{\mu}$) and smoothness ν as the function $\mathbf{x} \mapsto \sigma^2 K_{r,\nu}\left(\frac{\mathbf{x}}{\boldsymbol{\theta}}\right)$ (resp. $\mathbf{x} \mapsto \sigma^2 K_{r,\nu}(\mathbf{x}\boldsymbol{\mu})$).

Similarly, we define the Matérn tensorized covariance kernel with variance parameter σ^2 , correlation lengths $\boldsymbol{\theta}$ (resp. inverse correlation lengths $\boldsymbol{\mu}$) and smoothness ν as the function $\mathbf{x} \mapsto \sigma^2 K_{r,\nu}^{tens}\left(\frac{\mathbf{x}}{\boldsymbol{\theta}}\right)$ (resp. $\mathbf{x} \mapsto \sigma^2 K_{r,\nu}^{tens}(\mathbf{x}\boldsymbol{\mu})$).

Appendix 6.B Proofs of the existence of the Gibbs reference posterior

The proof of the existence and uniqueness of the Gibbs reference posterior that was used in Chapter 5 to deal with the Simple Kriging setting is inadequate in the Universal Kriging setting because the projection \mathbf{W}^\top may make key facts used in that chapter untrue. In the following, we provide replacements for the parts of the proof in Appendix 6.A that are invalid in the Universal Kriging setting.

The proof contained two parts, one dealing with “low correlations”, that is $\|\boldsymbol{\mu}\| \rightarrow +\infty$ and one with “high correlation”, that is $\|\boldsymbol{\mu}\| \rightarrow 0$.

Accounting for low correlation: $\|\boldsymbol{\mu}\| \rightarrow \infty$

Concerning the part about $\|\boldsymbol{\mu}\| \rightarrow +\infty$, we need to make sure that Corollary 5.13 remains true.

Define the functions h_i by

$$h_i(\mu_i | \boldsymbol{\mu}_{-i}) := \sqrt{\text{Tr} \left[\left(\frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}} \right)^2 \right]} = \left\| \frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}} \right\|. \quad (6.31)$$

The conclusion of Corollary 5.13 is that there exist $S > 0$ and $0 < a < b$ such that, whenever $\|\boldsymbol{\mu}\| \geq S$,

$$a h_i(\mu_i | \boldsymbol{\mu}_{-i}) \leq f_i(\mu_i | \boldsymbol{\mu}_{-i}) \leq b h_i(\mu_i | \boldsymbol{\mu}_{-i}). \quad (6.32)$$

We need to find conditions under which this is true. While the right inequality is obvious, the left inequality is harder to show.

Fix $\boldsymbol{\alpha} = \boldsymbol{\mu} / \|\boldsymbol{\mu}\|_{\infty}$. Then define $\mathbf{L}_{i,\boldsymbol{\alpha}} = \lim_{\|\boldsymbol{\mu}\| \rightarrow \infty} \frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}} / \left\| \frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_{\boldsymbol{\mu}} \right\|_{\infty}$.

We now give an explicit form for $\mathbf{L}_{i,\boldsymbol{\alpha}}$. Let \mathbf{X} be the $n \times r$ matrix representing the design set, and let $\mathbf{X}_{\boldsymbol{\alpha}}$ be the matrix $\mathbf{X} \text{Diag}(\boldsymbol{\alpha})$, where $\text{Diag}(\boldsymbol{\alpha})$ is the $r \times r$ diagonal matrix whose diagonal is the vector $\boldsymbol{\alpha}$.

Proposition 6.16. *If the Matérn kernel is anisotropic geometric, then $\mathbf{L}_{i,\boldsymbol{\alpha}}$ is the symmetric $n \times n$ matrix with null diagonal whose nondiagonal coefficients are given by the following rule: its (a, b) coefficient ($a, b \in \llbracket 1, n \rrbracket$ and $a \neq b$) is -1 if the a -th and b -th point in the design set $\mathbf{X}_{\boldsymbol{\alpha}}$ achieve minimal Euclidean distance within this design set, and 0 otherwise.*

Proof. We only prove the result when $\nu > 1$, but the proof is very similar in the case where $0 < \nu \leq 1$.

Abramowitz and Stegun [1964] (formula 9.7.2.) yields an equivalent for the one-dimensional Matérn kernel when $t \rightarrow +\infty$:

$$K_{1,\nu}(t) \sim \frac{\sqrt{\pi/2}}{\Gamma(\nu)2^{\nu-1}} (2\sqrt{\nu}t)^{\nu-1/2} \exp(-2\sqrt{\nu}t) \quad (6.33)$$

From Abramowitz and Stegun [1964] (formula 9.6.28.), we obtain that:

$$\begin{aligned} K'_{1,\nu}(t) &= -\frac{2\nu t}{\nu-1} K_{1,\nu-1} \left(\sqrt{\frac{\nu}{\nu-1}} t \right) \\ &\sim -2\sqrt{\nu} \frac{\sqrt{\pi/2}}{\Gamma(\nu)2^{\nu-1}} (2\sqrt{\nu}t)^{\nu-1/2} \exp(-2\sqrt{\nu}t) \\ &\sim -2\sqrt{\nu} K_{1,\nu}(t) \end{aligned} \quad (6.34)$$

The result follows after recalling that $\frac{\partial}{\partial \mu_i} K_{r,\nu}(\boldsymbol{\mu}\mathbf{x}) = \mu_i x_i^2 \|\boldsymbol{\mu}\mathbf{x}\|^{-1} K'_{1,\nu}(\|\boldsymbol{\mu}\mathbf{x}\|)$. When $\|\boldsymbol{\mu}\| \rightarrow \infty$,

$$\frac{\partial}{\partial \mu_i} K_{r,\nu}(\boldsymbol{\mu}\mathbf{x}) \sim -2\sqrt{\nu} \mu_i x_i^2 \|\boldsymbol{\mu}\mathbf{x}\|^{-1} \frac{\sqrt{\pi/2}}{\Gamma(\nu)2^{\nu-1}} (2\sqrt{\nu}\|\boldsymbol{\mu}\mathbf{x}\|)^{\nu-1/2} \exp(-2\sqrt{\nu}\|\boldsymbol{\mu}\mathbf{x}\|). \quad (6.35)$$

In the case where $0 < \nu \leq 1$, $\frac{\partial}{\partial \mu_i} K_{r,\nu}(\boldsymbol{\mu}\mathbf{x})$ also has an equivalent when $\|\boldsymbol{\mu}\| \rightarrow \infty$ whose prominent factor is $\exp(-2\sqrt{\nu}\|\boldsymbol{\mu}\mathbf{x}\|)$, so the end result is the same. \square

Proposition 6.17. *If the Matérn kernel is tensorized with smoothness $\nu > 1$, then $L_{i,\alpha}$ is the matrix with nonpositive coefficients such that $\|L_{i,\alpha}\|_\infty = 1$ which is proportional to the symmetric matrix described hereafter: it has null diagonal and its nondiagonal coefficients are given by the following rule: its (a, b) coefficient ($a, b \in \llbracket 1, n \rrbracket$ and $a \neq b$) is 0 if the a -th and b -th point in the design set \mathbf{X}_α do not achieve minimal 1-distance within this design set, and $\alpha_i^{\nu-1/2} |x_i^{(a)} - x_i^{(b)}|^{\nu+1/2} \prod_{j \neq i} \alpha_j^{\nu-1/2} |x_j^{(a)} - x_j^{(b)}|^{\nu-1/2}$ if they do.*

Remark. If the Matérn kernel is tensorized with smoothness $0 < \nu \leq 1$, then the same rule applies but with different formula when minimal 1-distance is achieved.

Proof. The proof is similar to that of Proposition 6.16. \square

Corollary 6.18. *For Matérn anisotropic geometric and tensorized kernels, if the design set \mathbf{X} is randomly chosen according to the Uniform probability distribution on $(0, 1)^{rn}$, then almost surely, whatever $i \in \llbracket 1, n \rrbracket$ and α in \mathbb{R}^r such that $\|\alpha\|_\infty = 1$, $L_{i,\alpha}$ has rank lower or equal to $2r$.*

Proof. Almost surely, whatever α in \mathbb{R}^r such that $\|\alpha\|_\infty = 1$, the design set \mathbf{X}_α has at most r couples of distinct points achieving equal distance (whether that distance be the 1- or 2-distance). *A fortiori*, it has at most r couples of distinct points achieving minimal distance. \square

With fixed α , as $\|\mu\| \rightarrow \infty$, we have

$$f_i(\mu_i | \mu_{-i}) \sim \left\| \frac{\partial}{\partial \mu_i} \Sigma_\mu \right\|_\infty \sqrt{\text{Tr} \left[\left(\mathbf{W}^\top L_{i,\alpha} \mathbf{W} \right)^2 \right] - \frac{1}{n-p} \text{Tr} \left[\mathbf{W}^\top L_{i,\alpha} \mathbf{W} \right]^2} \quad (6.36)$$

We may recognize the factor under the square root as the variance (multiplied by $n-p$) of the eigenvalues (accounting for multiplicity) of the matrix $\mathbf{W}^\top L_{i,\alpha} \mathbf{W}$. If the premise of Corollary 6.18 holds, and if $2r < n-p$, then it is null if and only if $\mathbf{W}^\top L_{i,\alpha} \mathbf{W}$ is the null matrix. Assumption 3 is designed to prevent this from happening.

Proposition 6.19. *Assume $2r < n-p$. For Matérn anisotropic geometric or tensorized correlation kernels, if the design set \mathbf{X} is randomly chosen according to the Uniform probability distribution on $(0, 1)^{rn}$, then almost surely, Assumption 3 implies that*

$$\min_{i \in \llbracket 1, n \rrbracket, \|\alpha\|_\infty = 1} \sqrt{\text{Tr} \left[\left(\mathbf{W}^\top L_{i,\alpha} \mathbf{W} \right)^2 \right] - \frac{1}{n-p} \text{Tr} \left[\mathbf{W}^\top L_{i,\alpha} \mathbf{W} \right]^2} > 0. \quad (6.37)$$

Proof. First, set $i \in \llbracket 1, n \rrbracket$ and α in \mathbb{R}^r such that $\|\alpha\|_\infty = 1$. We prove that $\mathbf{W}^\top L_{i,\alpha} \mathbf{W}$ is not the null matrix.

Assume that it is and that Assumption 3 holds. Assumption 3 implies that the intersection of the vector space spanned by \mathbf{P} and the image of $L_{i,\alpha}$ is $\{\mathbf{0}_n\}$. Therefore, for any $\mathbf{z} \in \mathbb{R}^{n-p}$, if $L_{i,\alpha} \mathbf{W} \mathbf{z} \neq \mathbf{0}_n$, then $\mathbf{W}^\top L_{i,\alpha} \mathbf{W} \mathbf{z} \neq \mathbf{0}_{n-p}$, which contradicts the assumption that $\mathbf{W}^\top L_{i,\alpha} \mathbf{W}$ is the null matrix. So $L_{i,\alpha} \mathbf{W}$ is the null $n \times (n-p)$ matrix, and thus the vector space spanned by \mathbf{W} is included in the kernel of $L_{i,\alpha}$. This implies that $L_{i,\alpha} \mathbf{P} \mathbf{P}^\top = L_{i,\alpha}$, and then that $\mathbf{P} \mathbf{P}^\top L_{i,\alpha} = L_{i,\alpha}$. However, per Propositions 6.16 and 6.17, all vectors in the image of $L_{i,\alpha}$ have at most $2r$ non-null elements when expressed in the canonical base of \mathbb{R}^n , so Assumption

3 implies that $\mathbf{P}^\top \mathbf{L}_{i,\alpha}$ is the null $p \times n$ matrix, and thus that $\mathbf{L}_{i,\alpha}$ is the null $n \times n$ matrix, which is untrue.

So, under Assumption 3, whatever $i \in \llbracket 1, n \rrbracket$ and α in \mathbb{R}^r such that $\|\alpha\|_\infty = 1$, $\mathbf{W}^\top \mathbf{L}_{i,\alpha} \mathbf{W}$ is not the null matrix and thus has a non-null eigenvalue. Moreover, $2r < n - p$ implies, according to Corollary 6.18, that it also almost surely has a null eigenvalue, so the standard deviation of its eigenvalues is positive. As the number of possible matrices $\mathbf{L}_{i,\alpha}$ (with $i \in \llbracket 1, n \rrbracket$ and α in \mathbb{R}^r such that $\|\alpha\|_\infty = 1$) is almost surely finite, this yields the result. \square

Corollary 6.20. *Assume $2r < n - p$. For Matérn anisotropic geometric or tensorized correlation kernels, if the design set \mathbf{X} is randomly chosen according to the Uniform probability distribution on $(0, 1)^{rn}$, then almost surely, Assumption 3 implies that there exist $S > 0$ and $0 < a < b$ such that whenever $\|\boldsymbol{\mu}\| \geq S$, Equation (6.32) holds.*

Accounting for high correlation: $\|\boldsymbol{\mu}\| \rightarrow 0$

In the part of the proof in Appendix 5.A concerning $\|\boldsymbol{\mu}\| \rightarrow 0$, we used a the series expansion of $\boldsymbol{\Sigma}_\boldsymbol{\mu}$. This expansion may be heavily modified by premultiplication by \mathbf{W}^\top and postmultiplication by \mathbf{W} .

In the case where $\nu < 1$, there is no material change unless the vector $\mathbf{1}$ belongs to the vector space spanned by \mathbf{H} .

Proof of Proposition 6.8. Because $\mathbf{1}$ does not belong to the vector space spanned by \mathbf{H} , $\mathbf{W}^\top \mathbf{1} \mathbf{1}^\top \mathbf{W}$ has rank 1 and so the proof of this result is the same as in the Simple Kriging case. \square

If $\mathbf{1}$ *does* belong to the vector space spanned by \mathbf{H} , further study would be needed to assess whether or not the above theorem still applies, essentially because we cannot count on $L(\mathbf{y}|\boldsymbol{\mu})$ vanishing as $\|\boldsymbol{\mu}\| \rightarrow 0$.

Let us now focus on the case where $\nu > 1$. We reproduce in Lemma 6.21 key facts given by Lemma 5.7 and Proposition 5.22:

Lemma 6.21. *For any Matérn anisotropic geometric or tensorized correlation kernel with smoothness parameter $\nu > 1$, if a coordinate-distinct design set is used, there exists a > 0 such that when $\|\boldsymbol{\mu}\| \rightarrow 0$:*

1. $\left\| \frac{\partial}{\partial \mu_i} \boldsymbol{\Sigma}_\boldsymbol{\mu} \right\| = O(\mu_i)$;
2. $\left\| \boldsymbol{\Sigma}_\boldsymbol{\mu}^{-1} \right\| = O(\|\boldsymbol{\mu}\|^{-a})$.

The newt result follows immediately.

Corollary 6.22. *For any Matérn anisotropic geometric or tensorized correlation kernel with smoothness parameter $\nu > 1$, if a coordinate-distinct design set is used, there exist a > 0 , $m > 0$ and $S > 0$ such that, for any $\boldsymbol{\mu} \in (0, +\infty)^r$ such that $\|\boldsymbol{\mu}\| \leq S$ and $\mu_i \leq \|\boldsymbol{\mu}\|^a$, $f_i(\mu_i|\boldsymbol{\mu}_{-i}) \leq m$.*

We combine the previous fact with a useful universal majoration of $f_i(\mu_i|\boldsymbol{\mu}_{-i})$.

Proposition 6.23. *For an r -dimensional anisotropic geometric or tensorized Matérn correlation kernel with smoothness parameter ν pertaining to a design set containing n coordinate-distinct points, $\forall \boldsymbol{\mu} \in [0, +\infty)^r$ such that $\mu_i > 0$,*

$$f_i(\mu_i | \boldsymbol{\mu}_{-i}) \leq (n-p)(2\nu+r)\mu_i^{-1} \quad (6.38)$$

Proof. Whatever $x, y \in \mathbb{R}$, $K_{1,\nu}(x-y) = \int_{\mathbb{R}} \widehat{K}_{1,\nu}(\omega) e^{i\omega(x-y)} d\omega$.

For the sake of concision, we only consider the case where the Matérn kernel is anisotropic geometric, as the changes in the case of a tensorized kernel are straightforward.

Moreover, we start by proving the result in the case where \mathbf{W} is the identity matrix \mathbf{I}_n (Simple Kriging case).

$$\begin{aligned} \sum_{j,k=1}^n \xi_j \xi_k K_{r,\nu} \left((\mathbf{x}^{(j)} - \mathbf{x}^{(k)}) \boldsymbol{\mu} \right) &= \int_{\mathbb{R}^r} \widehat{K}_{r,\nu}(\boldsymbol{\omega}) \left| \sum_{j=1}^n \xi_j e^{i\boldsymbol{\omega}_i \mu_i x_i^{(j)} + i \langle \boldsymbol{\omega}_{-i} | \boldsymbol{\mu}_{-i} \mathbf{x}^{(j)} \rangle} \right|^2 d\boldsymbol{\omega} \\ &= M_r(\nu) \mu_i^{-1} I_{\boldsymbol{\mu}}(\boldsymbol{\xi}), \end{aligned} \quad (6.39)$$

where

$$M_r(\nu) = \frac{\Gamma(\nu + \frac{r}{2})(2\sqrt{\nu})^{2\nu}}{\pi^{\frac{r}{2}} \Gamma(\nu)}; \quad (6.40)$$

$$I_{\boldsymbol{\mu}}(\boldsymbol{\xi}) = \int_{\mathbb{R}^r} \left(4\nu + \mu_i^{-2} s_i^2 + \left\| \frac{\mathbf{s}_{-i}}{\boldsymbol{\mu}_{-i}} \right\|^2 \right)^{-\frac{r}{2}-\nu} \left| \sum_{j=1}^n \xi_j e^{i \langle \mathbf{s} | \mathbf{x}^{(j)} \rangle} \right|^2 d\mathbf{s}. \quad (6.41)$$

We also have

$$\frac{d}{d\mu_i} \sum_{j,k=1}^n \xi_j \xi_k K_{r,\nu} \left((\mathbf{x}^{(j)} - \mathbf{x}^{(k)}) \boldsymbol{\mu} \right) = -M_r(\nu) \mu_i^{-1} I_{\boldsymbol{\mu}}(\boldsymbol{\xi}) + M_r(\nu) \mu_i^{-1} \frac{d}{d\mu_i} I_{\boldsymbol{\mu}}(\boldsymbol{\xi}) \quad (6.42)$$

$$\begin{aligned} &\frac{d}{d\mu_i} I_{\boldsymbol{\mu}}(\boldsymbol{\xi}) \\ &= 2 \left(\frac{r}{2} + \nu \right) \mu_i^{-3} \int_{\mathbb{R}^r} s_i^2 \left(4\nu + \mu_i^{-2} s_i^2 + \left\| \frac{\mathbf{s}_{-i}}{\boldsymbol{\mu}_{-i}} \right\|^2 \right)^{-\frac{r}{2}-\nu-1} \left| \sum_{j=1}^n \xi_j e^{i \langle \mathbf{s} | \mathbf{x}^{(j)} \rangle} \right|^2 d\mathbf{s} \\ &= (2\nu+r) \mu_i^{-3} \int_{\mathbb{R}^r} \frac{s_i^2}{4\nu + \mu_i^{-2} s_i^2 + \left\| \frac{\mathbf{s}_{-i}}{\boldsymbol{\mu}_{-i}} \right\|^2} \left(4\nu + \mu_i^{-2} s_i^2 + \left\| \frac{\mathbf{s}_{-i}}{\boldsymbol{\mu}_{-i}} \right\|^2 \right)^{-\frac{r}{2}-\nu} \left| \sum_{j=1}^n \xi_j e^{i \langle \mathbf{s} | \mathbf{x}^{(j)} \rangle} \right|^2 d\mathbf{s} \end{aligned} \quad (6.43)$$

From this, we obtain that for any non-null vector $\boldsymbol{\xi} \in \mathbb{R}^n$,

$$0 < \frac{d}{d\mu_i} I_{\boldsymbol{\mu}}(\boldsymbol{\xi}) \leq (2\nu+r) \mu_i^{-1} I_{\boldsymbol{\mu}}(\boldsymbol{\xi}) \quad (6.44)$$

Now let us define the matrix $\mathbf{F}_{\boldsymbol{\mu}}$ as the matrix representing in the canonical base of \mathbb{R}^n the positive definite quadratic form $\boldsymbol{\xi} \mapsto M_r(\nu) \mu_i^{-1} \frac{d}{d\mu_i} I_{\boldsymbol{\mu}}(\boldsymbol{\xi})$. From the previous calculations, we

gather that $\frac{d}{d\mu_i}\Sigma_\mu = -\mu_i^{-1}\Sigma_\mu + \mathbf{F}_\mu$. This in turn yields $\left(\frac{\partial}{\partial\mu_i}\Sigma_\mu\right)\Sigma_\mu^{-1} = -\mu_i^{-1}\mathbf{I}_n + \mathbf{F}_\mu\Sigma_\mu^{-1}$ and $\left(\left(\frac{\partial}{\partial\mu_i}\Sigma_\mu\right)\Sigma_\mu^{-1}\right)^2 = \mu_i^{-2}\mathbf{I}_n + (\mathbf{F}_\mu\Sigma_\mu^{-1})^2 - 2\mu_i^{-1}\mathbf{F}_\mu\Sigma_\mu^{-1}$.

$$\mathrm{Tr}\left[\left(\frac{\partial}{\partial\mu_i}\Sigma_\mu\right)\Sigma_\mu^{-1}\right] = -n\mu_i^{-1} + \mathrm{Tr}[\mathbf{F}_\mu\Sigma_\mu^{-1}]. \quad (6.45)$$

$$\mathrm{Tr}\left[\left(\left(\frac{\partial}{\partial\mu_i}\Sigma_\mu\right)\Sigma_\mu^{-1}\right)^2\right] = n\mu_i^{-2} + \mathrm{Tr}[(\mathbf{F}_\mu\Sigma_\mu^{-1})^2] - 2\mu_i^{-1}\mathrm{Tr}[\mathbf{F}_\mu\Sigma_\mu^{-1}]. \quad (6.46)$$

$$\mathrm{Tr}\left[\left(\left(\frac{\partial}{\partial\mu_i}\Sigma_\mu\right)\Sigma_\mu^{-1}\right)^2\right] - \frac{1}{n}\mathrm{Tr}\left[\left(\frac{\partial}{\partial\mu_i}\Sigma_\mu\right)\Sigma_\mu^{-1}\right]^2 = \mathrm{Tr}[(\mathbf{F}_\mu\Sigma_\mu^{-1})^2] - \frac{1}{n}\mathrm{Tr}[\mathbf{F}_\mu\Sigma_\mu^{-1}]^2. \quad (6.47)$$

\mathbf{F}_μ and Σ_μ^{-1} being two symmetric positive definite matrices, their product $\mathbf{F}_\mu\Sigma_\mu^{-1}$ is diagonalizable and all its eigenvalues are positive. Thus $\mathrm{Tr}[(\mathbf{F}_\mu\Sigma_\mu^{-1})^2] \leq \mathrm{Tr}[\mathbf{F}_\mu\Sigma_\mu^{-1}]^2$.

Let $(\xi_\mu^j)_{1 \leq j \leq n}$ be a basis of unit eigenvectors of Σ_μ^{-1} . Then

$$\mathrm{Tr}[\mathbf{F}_\mu\Sigma_\mu^{-1}] = \sum_{j=1}^n (\xi_\mu^j)^\top \mathbf{F}_\mu\Sigma_\mu^{-1}\xi_\mu^j = \sum_{j=1}^n \frac{(\xi_\mu^j)^\top \mathbf{F}_\mu\xi_\mu^j}{(\xi_\mu^j)^\top \Sigma_\mu\xi_\mu^j} \leq n(2\nu + r)\mu_i^{-1}. \quad (6.48)$$

This implies that

$$\mathrm{Tr}\left[\left(\left(\frac{\partial}{\partial\mu_i}\Sigma_\mu\right)\Sigma_\mu^{-1}\right)^2\right] - \frac{1}{n}\mathrm{Tr}\left[\left(\frac{\partial}{\partial\mu_i}\Sigma_\mu\right)\Sigma_\mu^{-1}\right]^2 \leq n(n-1)(2\nu+r)^2\mu_i^{-2} \quad (6.49)$$

$$\sqrt{\mathrm{Tr}\left[\left(\left(\frac{\partial}{\partial\mu_i}\Sigma_\mu\right)\Sigma_\mu^{-1}\right)^2\right] - \frac{1}{n}\mathrm{Tr}\left[\left(\frac{\partial}{\partial\mu_i}\Sigma_\mu\right)\Sigma_\mu^{-1}\right]^2} \leq n(2\nu+r)\mu_i^{-1} \quad (6.50)$$

Now, if \mathbf{W} is not the identity matrix, then the previous proof still holds, albeit with some alterations. Instead of considering all non-null vectors $\xi \in \mathbb{R}^n$, we consider only those which can be expressed as $\mathbf{W}\xi_{\mathbf{W}}$, with $\xi_{\mathbf{W}}$ belonging to \mathbb{R}^{n-p} . In the same vein, once it comes to computing $\mathrm{Tr}\left[\mathbf{W}^\top \mathbf{F}_\mu \mathbf{W} (\mathbf{W}^\top \Sigma_\mu \mathbf{W})^{-1}\right]$, we use a basis $(\xi_{\mathbf{W},\mu}^j)_{1 \leq j \leq n-p}$ of unit eigenvectors of $(\mathbf{W}^\top \Sigma_\mu \mathbf{W})^{-1}$.

□

Proposition 6.24. *With a Matérn anisotropic geometric or tensorized correlation kernel with smoothness $\nu > 1$, if a design set with coordinate-distinct points is used, then Assumption 4 implies that there exists $\epsilon' > 0$ such that $L(\mathbf{y}|\boldsymbol{\mu})f_i(\mu_i|\boldsymbol{\mu}_{-i}) = O(\mu_i^{-1+\epsilon'})$ when $\|\boldsymbol{\mu}\| \rightarrow 0$.*

Proof. Assumption 4 ensures that $L(\mathbf{y}|\boldsymbol{\mu})$ is bounded as a function of $\boldsymbol{\mu}$. Because of Corollary 6.22, and using said Corollary's notations, we know that there exists $M > 0$ such that, for any $\boldsymbol{\mu} \in (0, +\infty)^r$ such that $\mu_i \leq \|\boldsymbol{\mu}\|^\alpha$, $L(\mathbf{y}|\boldsymbol{\mu})f_i(\mu_i|\boldsymbol{\mu}_{-i}) \leq M$.

Let us now focus on the $\boldsymbol{\mu} \in (0, +\infty)^r$ such that $\mu_i \geq \|\boldsymbol{\mu}\|^\alpha$. Then $\|\boldsymbol{\mu}\|^\epsilon \leq \mu_i^{\epsilon/a}$. Choosing $\epsilon' = \epsilon/a$, combining Assumption 4 and Proposition 6.23 yields the result. □

With the help of Proposition 6.24, using essentially the proof of Proposition 5.26, we obtain the following result.

Proposition 6.25. *In a Universal Kriging model with a Matérn anisotropic geometric or tensorized correlation kernel with smoothness $\nu > 1$, if a design set with coordinate-distinct points is used, then Assumption 4 implies that the conditional posterior distribution $\pi_i(\mu_i | \mathbf{y}, \boldsymbol{\mu}_{-i})$, seen as a function of $\boldsymbol{\mu}$, is continuous over $\{\boldsymbol{\mu} \in [0, +\infty)^r : \mu_i \neq 0\}$.*

Proofs of the main results

Proof of Theorem 6.9. With the help of Proposition 6.25 and Corollary 6.20, a similar result to Lemma 5.27 can be proved. Then the proof of Theorem 6.9 works like the proof of Theorem 5.2. \square

In order to be able to state the remaining results in a concise manner, we implicitly assume from this point onwards that the correlation kernel is Matérn anisotropic geometric or tensorized. Consider the following set of conditions:

1. $1 < \nu < 2$ and $n > r + 1$;
2. $2 < \nu < 3$ and $n > (r + 1)(r/2 + 2)$.

Now define k_ν as in the previous chapter: as the orthogonal complement in \mathbb{R}^n of the vector space spanned by:

1. if $\nu \in (1, 2)$: $\mathbf{1}$ and \mathbf{X}_i ($i \in \llbracket 1, r \rrbracket$);
2. if $\nu \in (2, 3)$: $\mathbf{1}$ and \mathbf{X}_i ($i \in \llbracket 1, r \rrbracket$) and $\mathbf{X}_i \circ \mathbf{X}_j$ ($i, j \in \llbracket 1, r \rrbracket$).

Proposition 6.26. *In the case of Universal Kriging where the mean function space is included within the space of polynomials of degree 0 or 1, if one of the previous conditions is satisfied, for any vector $\mathbf{y} \in \mathbb{R}^n$ not orthogonal to k_ν , when $\|\boldsymbol{\mu}\| \rightarrow 0$,*

$$\|\boldsymbol{\mu}\|^{-2\nu} = O\left(\mathbf{y}^\top \mathbf{W} \left(\mathbf{W}^\top \boldsymbol{\Sigma}_\mu \mathbf{W}\right)^{-1} \mathbf{W}^\top \mathbf{y}\right). \quad (6.51)$$

Proof. The proof is similar to that of Proposition 5.23 after noticing that $\mathbf{W}\mathbf{W}^\top$ is a projection on a vector space that contains k_ν . \square

Proposition 6.27. *In the case of Ordinary Kriging, under the conditions of Theorem 6.9, if one of the previous conditions is satisfied, then there exists a hyperplane \mathcal{H} of \mathbb{R}^n such that, provided $\mathbf{y} \in \mathbb{R}^n \setminus \mathcal{H}$, Assumption 4 is true.*

Proposition 6.28. *In the case of Universal Kriging where the mean function space is included within the space of polynomials of degree 0 or 1, if $2 < \nu < 3$ and $n > r(r + 1)/2 + 2r + 3$, then there exists a hyperplane \mathcal{H} of \mathbb{R}^n such that, provided $\mathbf{y} \in \mathbb{R}^n \setminus \mathcal{H}$, Assumption 4 is true.*

Proof of Propositions 6.27 and 6.28. Let $v_1(\boldsymbol{\mu}) \geq v_2(\boldsymbol{\mu}) \geq \dots \geq v_{n-p}(\boldsymbol{\mu})$ be the ordered eigenvalues of $\mathbf{W}^\top \boldsymbol{\Sigma}_\mu \mathbf{W}$. We can now rewrite $L(\mathbf{y} | \boldsymbol{\mu})$ as

$$L(\mathbf{y} | \boldsymbol{\mu})^2 \propto \prod_{k=1}^{n-p} \left[v_k(\boldsymbol{\mu})^{-1} \left(\mathbf{y}^\top \mathbf{W} \left(\mathbf{W}^\top \boldsymbol{\Sigma}_\mu \mathbf{W}\right)^{-1} \mathbf{W}^\top \mathbf{y} \right)^{-1} \right]. \quad (6.52)$$

Proposition 6.26 asserts that for any $\mathbf{y} \in \mathbb{R}^n$ that is not orthogonal to k_ν , the following holds: $\left(\mathbf{y}^\top \mathbf{W} \left(\mathbf{W}^\top \boldsymbol{\Sigma}_\mu \mathbf{W}\right)^{-1} \mathbf{W}^\top \mathbf{y} \right)^{-1} = O(\|\boldsymbol{\mu}\|^{2\nu})$ for $\|\boldsymbol{\mu}\| \rightarrow 0$.

Besides, Proposition 5.22 asserts that $\|\Sigma_{\boldsymbol{\mu}}^{-1}\| = O(\|\boldsymbol{\mu}\|^{-2\nu})$, so a fortiori $\left\| \left(\mathbf{W}^\top \Sigma_{\boldsymbol{\mu}} \mathbf{W} \right)^{-1} \right\| = O(\|\boldsymbol{\mu}\|^{-2\nu})$ and therefore $\left(\mathbf{y}^\top \mathbf{W} \left(\mathbf{W}^\top \Sigma_{\boldsymbol{\mu}} \mathbf{W} \right)^{-1} \mathbf{W}^\top \mathbf{y} \right)^{-1} = O \left(\left\| \left(\mathbf{W}^\top \Sigma_{\boldsymbol{\mu}} \mathbf{W} \right)^{-1} \right\|^{-1} \right)$.

This implies that for every integer $i \in \llbracket 1, r \rrbracket$, $v_k(\boldsymbol{\mu})^{-1} \left(\mathbf{y}^\top \mathbf{W} \left(\mathbf{W}^\top \Sigma_{\boldsymbol{\mu}} \mathbf{W} \right)^{-1} \mathbf{W}^\top \mathbf{y} \right)^{-1} = O(1)$.

Moreover, $\left\| \mathbf{W}^\top \Sigma_{\boldsymbol{\mu}} \mathbf{W} \right\| = O(v_1(\boldsymbol{\mu}))$ so $v_1(\boldsymbol{\mu})^{-1} = O \left(\left\| \mathbf{W}^\top \Sigma_{\boldsymbol{\mu}} \mathbf{W} \right\|^{-1} \right)$. In all cases considered, $\|\boldsymbol{\mu}\|^{2\lfloor \nu \rfloor} = O \left(\left\| \mathbf{W}^\top \Sigma_{\boldsymbol{\mu}} \mathbf{W} \right\| \right)$, which implies that the following asymptotic upper bound holds: $v_1(\boldsymbol{\mu})^{-1} \left(\mathbf{y}^\top \mathbf{W} \left(\mathbf{W}^\top \Sigma_{\boldsymbol{\mu}} \mathbf{W} \right)^{-1} \mathbf{W}^\top \mathbf{y} \right)^{-1} = O \left(\|\boldsymbol{\mu}\|^{2\nu - 2\lfloor \nu \rfloor} \right)$.

In the end,

$$L(\mathbf{y}|\boldsymbol{\mu})^2 = O \left(\|\boldsymbol{\mu}\|^{2\nu - 2\lfloor \nu \rfloor} \right). \quad (6.53)$$

□

Proof of Propositions 6.10 and 6.11. Propositions 6.10 and 6.11 are obtained by combining Theorem 6.9 with Propositions 6.27 and 6.28 respectively. □

Chapter 7

Trans-Gaussian Kriging in a Bayesian framework: a case study

This chapter corresponds to the article Muré [2018b].

Abstract

In the context of Gaussian Process Regression or Kriging, we propose a full-Bayesian solution to deal with hyperparameters of the covariance function. This solution can be extended to the Trans-Gaussian Kriging framework, which makes it possible to deal with spatial data sets that violate assumptions required for Kriging. It is shown to be both elegant and efficient. We propose an application to computer experiments in the field of nuclear safety, where it is necessary to model non-destructive testing procedures based on eddy currents to detect possible wear in steam generator tubes.

Résumé

Dans le but de faciliter la régression par processus gaussiens ou krigeage, nous proposons une solution pleinement bayésienne pour traiter les hyperparamètres de la fonction de covariance. Cette solution peut être étendue au krigeage trans-gaussien, ce qui permet d'utiliser des jeux de données qui violent les hypothèses du krigeage. La méthode se révèle à la fois élégante et efficace. Nous proposons une application aux simulations numériques dans le domaine de la sûreté nucléaire, domaine dans lequel il est nécessaire de modéliser les tests non destructifs fondés sur les courants de Foucault visant à détecter de possibles usures dans les tubes de générateurs de vapeur.

7.1 Introduction

Non-destructive testing (NDT) is a group of techniques used in industry to evaluate the properties of components or systems without destroying them. Numerical simulations can be used to calibrate NDT techniques and determine adequate threshold levels for defect detection. These simulations can be computationally expensive. As a result, we may want to replace them by a cheaper surrogate model.

The NDT problem that motivates this chapter is the following one. In order to improve nuclear safety, Électricité de France (EDF) has developed a Non-Destructive Evaluation for testing the presence of a defect in Steam Generator Tubes. In nuclear power plants, steam generators

are the interface between primary and secondary water circuits. High-pressured water from the primary circuit flows into the steam generator tubes. To prevent the tubes from vibrating under this solicitation, they are held in place by anti-vibration bars (AVB). Rubbing may in time leave defects in the tubes, however. To detect them, a SAX (axial) probe is moved down the tube. This Non-Destructive Exam was modeled in C3D, a computer code which can accurately represent any possible defect geometry [Maurice et al., 2013]. As this code uses the finite element method, a wide range of parameters may be taken into account. Furthermore, any degree of accuracy can be reached as long as the mesh is fine enough [Thomas et al., 2015]. Non-destructive inspections based on eddy currents exploit the way the induction flux changes as the probe approaches a defect. If the tube were a perfect cylinder, the coils of the probe would get the same flux by cylindrical symmetry. A large enough difference between both fluxes therefore signals a defect. This differential flux is a complex quantity, but in the presented application only its imaginary part is used for defect detection, or rather the difference between maximum and minimum imaginary part as the probe moves through the tube. It is then converted to a tension by Lenz’s law.

Following both expert reports and data simulations, eight geometrical parameters (see Figure 7.1) and one non-geometrical parameter were identified as having the greatest influence on the output of C3D. In order to be able to define POD curves, they were given probability distributions reflecting expert opinion.

1. $a \sim U[a_a, b_a]$: defect depth (mm);
2. $E \sim N(a_e, b_e)$: pipe thickness (mm) based on measurements of 5000 pipes;
3. $h_1 \sim U[a_{h_1}, b_{h_1}]$: distance between the AVB and the top of the defect (mm);
4. $h_2 \sim U[a_{h_2}, b_{h_2}]$: distance between the AVB and the bottom of the defect (mm);
5. $e_{BAV1} \sim U[-a + a_{e_{BAV1}}, b_{e_{BAV1}}]$: length of the gap between the AVB on the side of the defect and the tube (mm);
6. $e_{BAV2} \sim U[a_{e_{BAV2}}, b_{e_{BAV2}}]$: length of the gap between the AVB on the other side and the tube (mm);
7. $h_{BAV} \sim U[a_{h_{BAV}}, b_{h_{BAV}}]$: shift between both AVBs (mm);
8. $inc \sim \mathcal{N}_{trunc,0}(inc_a, inc_b)$: inclination of the AVB on the side of the defect.
 $\mathcal{N}_{trunc,0}(inc_a, inc_b)$ denotes a Normal distribution with mean inc_a and variance inc_b truncated at 0: its support is $[0, +\infty)$ (mm);
9. $cond \sim \mathcal{N}(cond_a, cond_b)$: conductivity of the tube (MS/m).

The parameter of interest being a , all other parameters are collectively denoted by the eight-dimensional vector $\mathbf{x} = (E, h_1, h_2, e_{BAV1}, e_{BAV2}, h_{BAV}, inc, cond)^\top$. \mathbf{X} is the corresponding random variable: its joint distribution is the product of the distributions described above. For ease of manipulation, \mathbf{X} is reparametrized in such a way as to follow the Uniform distribution on the 8-dimensional cube $(0, 1)^8$. We also reparametrize a in the same way. The whole input space becomes the unit cube $(0, 1)^r$ with $r = 9$.

Our ultimate goal is to compute the Probability Of Detection (POD) of a defect as a function of its depth a . We formalize this notion in Section 7.2. Practical computation of the POD requires the use of surrogate models. Section 7.3 provides the mathematical framework for dealing with surrogate model uncertainty. Finally, Section 7.4 presents the resulting POD curves in the case of steam generator defect detection.

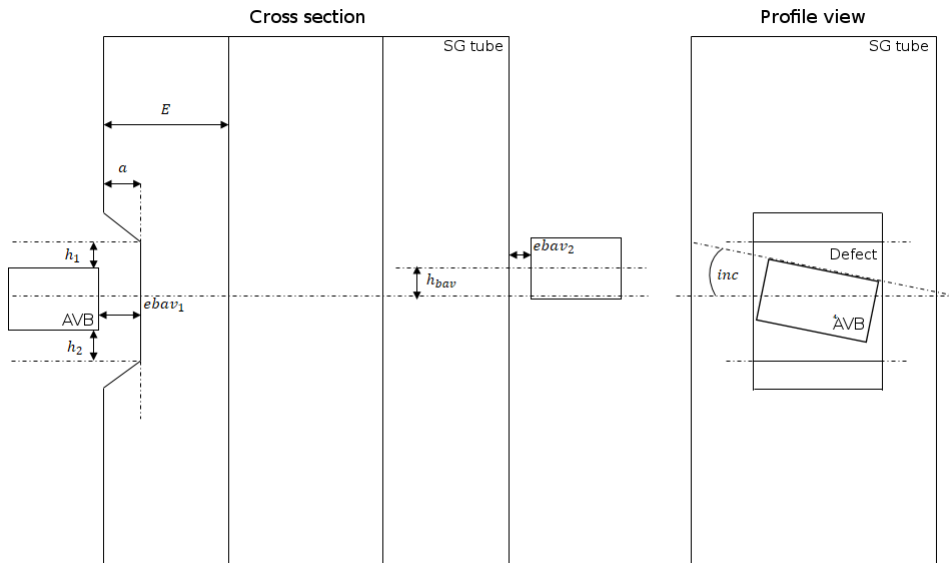


Figure 7.1 – Geometrical parameters of the computer model C3D. Left: Profile view of a steam generator tube where the left AVB caused wear. Right: View of the side with the defect.

7.2 Probability Of Detection (POD)

Industrial practice often uses Functional Risk Curves (FRC) as a means of expressing the “probability” of some (un)desirable event based on the value of some critical parameter. Probability Of Detection (POD) curves are a particular kind of FRC. They represent the probability that a given testing procedure has of detecting a defect as a function of some parameter of interest. Two factors can justify the probabilistic framework:

1. the testing procedure may incorporate some randomness;
2. the result may depend not only on the parameter of interest but also on unknown nuisance parameters.

We denote by a the parameter of interest and by \mathbf{x} the nuisance parameters. The set of all values possibly taken by \mathbf{x} is endowed with a probability distribution which should reflect their frequency in real life. Let \mathbf{X} be a random variable following this probability distribution. For given a and \mathbf{x} , $z(a, \mathbf{x})$ denotes the output of the testing procedure. Depending on the specific testing procedure used, it may or may not be random. In any case, it falls to the modeler to specify the probability distribution of \mathbf{X} and/or the probability distributions of $z(a, \mathbf{x})$ for all a and \mathbf{x} . In this chapter, we only deal with deterministic testing procedures so POD curves are entirely determined by the distribution of \mathbf{X} .

Typically, the output z is a scalar quantity. When this quantity lies beyond a predetermined threshold s , it signals the presence of a defect, so the POD curve is the function defined by

$$POD(a) = \mathbb{P}(z(a, \mathbf{X}) > s). \quad (7.1)$$

Being a deterministic mapping does not make z easy to determine, however. Computing the POD curve would in theory require conducting physical tests on a wide sample of objects with a wide variety of defects – i.e. a wide variety of parameters a and \mathbf{x} . Because of the

prohibitive costs involved, phenomenological numerical simulations are used to simulate the testing procedure. Throughout this chapter, such simulations – collectively called a “computer experiment” – are assumed to be perfectly accurate.

Unfortunately, even these numerical simulations are too costly to allow a Monte-Carlo approach. Therefore, \mathbf{z} needs to be approximated by some less costly model. In order to control the ensuing error, statistical surrogate models are used – linear regression or polynomial chaos for example. In this chapter, we focus on Kriging, otherwise known as Gaussian Process regression, and on Trans-Gaussian Kriging, where a transformation step is performed on the output space before Kriging is performed. When a surrogate model is used, z must be replaced by Z , a random mapping whose probability distribution represents the uncertainty about z .

Surrogate models introduce some ambiguity in the definition of the POD curve. Should it be approximated by $POD(a) \approx \mathbb{P}(Z(a, \mathbf{X}) > s)$? Although natural, this approximation conflates the uncertainty about \mathbf{X} and about Z , although the two are different in nature. The distribution of \mathbf{X} is understood in *frequentist* terms: over the year, an experimenter conducting tests on multiple pieces of equipment will encounter varied values of \mathbf{X} . The distribution of Z represents *epistemic* uncertainty, which is irrelevant to the actual tests conducted in real life and could be mitigated if more computational resources were available. This point is illustrated in the description of a concrete application below.

Kriging and Trans-Gaussian Kriging surrogate models depend on parameters, which will henceforth be named “hyperparameters” to differentiate them from the “physical” parameters a and \mathbf{X} . Hyperparameters are tricky, because they have tremendous influence so careful surrogate model calibration is normally required. The Bayesian paradigm, being a complete inferential approach (Robert et al. [2011] provides a review of its range from a practical standpoint) can be used to avoid this step. In this chapter, we propose eliminating hyperparameters from the model by means of Bayesian averaging out. Prior elicitation is done with the help of Bernardo’s reference prior theory [Bernardo, 2005].

7.3 An Objective Bayesian outlook to Trans-Gaussian Kriging

Kriging Likelihood

For the sake of simplicity, we assume that the parameter of interest a belongs to $[0, 1]$ and that \mathbf{x} , which represents all nuisance parameters belongs to $[0, 1]^{r-1}$ for some positive integer r . The scalar output of interest is $y(a, \mathbf{x})$.

Kriging is a very flexible surrogate model for computer experiments [Santner et al., 2003]. To use it, the computer experiment must first be conducted for $n \in \mathbb{N}$ values $(a^{(i)}, \mathbf{x}^{(i)})$ (called observations points) of the parameters. The set of observation points $((a^{(i)}, \mathbf{x}^{(i)}))_{i \in \llbracket 1, n \rrbracket}$ is called the design set. The vector of outputs $\mathbf{y} = (y(a^{(1)}, \mathbf{x}^{(1)}), \dots, y(a^{(n)}, \mathbf{x}^{(n)}))^\top \in \mathbb{R}^n$ is called the observation vector.

Kriging models the uncertainty about y by defining Y as a Gaussian process subjected to the condition that for every integer $i \in \llbracket 1, n \rrbracket$, $Y(a, \mathbf{x}_i) = y(a, \mathbf{x}_i)$. Because this conditioning is critical, we need a way to differentiate the conditioned from the unconditioned process, so let \mathcal{Y} be the unconditioned Gaussian process.

In the literature, \mathcal{Y} is often assumed to be stationary. That is, the distribution of the unconditioned Gaussian Process (before knowing the observed values) should be invariant by translation. This simplifying assumption is often reasonable, because should the emulated function be non-stationary, the distribution of the Gaussian Process conditional to the observations would reflect this non-stationarity, provided the number of observations is sufficient. Moreover, assuming non-stationarity would require some sort of prior knowledge of the kind of non-stationarity that is expected. For example, warped Gaussian processes [Snelson et al., 2004] use the knowledge of the presence of a discontinuity.

In some contexts, an additional assumption can be made: that the distribution of the Gaussian Process before observing the data is isotropic. When it is possible to make this assumption, the correlation kernel of the Gaussian Process is usually parametrized by two positive hyper-parameters: variance σ^2 and correlation length θ . Seeing the Gaussian Process as a response surface, one may think of $\sigma := \sqrt{\sigma^2}$ as the scale of variation of the output and of θ as the scale of variation of the input.

While isotropy is a natural assumption in geostatistics – the original application of Kriging [Matheron, 1960] – it is rarely adequate when dealing with computer experiments. Each of the r dimensions in the input space $[0, 1]^r$ corresponds to one parameter, and the parameters can be heterogeneous. In such contexts, a correlation length θ_i is necessary for every parameter ($i \in \llbracket 1, r \rrbracket$). The covariance kernel is then anisotropic and we denote by $\boldsymbol{\theta}$ the vector $(\theta_i)_{i=1}^r$. The best kind of anisotropic kernel for interpretation is anisotropic geometric, but tensorized kernels are often used for simplicity [Stein, 1999, page 54].

The easiest way to introduce non-stationarity is through the addition of some non-constant deterministic function. In the Universal Kriging framework, this function is assumed to belong to a given small-dimensional vectorial space \mathcal{F}_p . One assumes therefore that there exists in \mathcal{F}_p a function f that adequately approximates the deterministic function one seeks to emulate, and the stationary Gaussian process then merely models our perception of the error made when using such an approximation.

Let $(f_j)_{j=1}^p$ be a basis of \mathcal{F}_p and let $\boldsymbol{\beta} = (\beta_j)_{j=1}^p \in \mathbb{R}^p$ be the vector such that $f = \sum_{j=1}^p \beta_j f_j$. Naturally, p should be smaller than the number of observation points n or the model would not be identifiable. Denote by \mathbf{H} the $n \times p$ matrix whose (i, j) -th element is $f_j(\mathbf{x}^{(i)})$. Again for the sake of identifiability, assume that \mathbf{H} is of full rank. Let then \mathbf{y} be the vector of length n whose i -th element is the observation of the Gaussian process made at $\mathbf{x}^{(i)}$. It is a Gaussian vector with mean vector $\mathbf{H}\boldsymbol{\beta}$. Let $\boldsymbol{\Sigma}_\theta$ be its correlation matrix.

The likelihood of the parameters $\boldsymbol{\beta}$, σ^2 and $\boldsymbol{\theta}$ when observing \mathbf{y} is then:

$$L(\mathbf{y} \mid \boldsymbol{\beta}, \sigma^2, \boldsymbol{\theta}) = \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{n}{2}} |\boldsymbol{\Sigma}_\theta|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{H}\boldsymbol{\beta})^\top \boldsymbol{\Sigma}_\theta^{-1} (\mathbf{y} - \mathbf{H}\boldsymbol{\beta}) \right\}. \quad (7.2)$$

Transformation

If even more flexibility is needed, like in some cases arising naturally from examples in the natural sciences – see the example below – one can relax the stationarity assumption and even the assumption that the random process is Gaussian through Trans-Gaussian Kriging [De Oliveira et al., 1997]. The idea is to assume that a random field $\mathcal{Z}(a, \mathbf{x})$ would be Gaussian

and stationary if the output space O (which is assumed to be an interval of \mathbb{R}) were transformed in an adequate fashion. Practically speaking, one must choose a parametric family of nondecreasing differentiable transformations $g_\alpha : O \rightarrow \mathbb{R}$. For some α , $\mathcal{Y} := g_\alpha(\mathcal{Z})$ is a stationary Gaussian random field.

The detection of defects on Steam Generator tubes is most importantly impacted by the depth of the defect. As a first approximation, one may say that the greater the defect depth, the greater the chances of detecting it.

Figure 7.2 (left) illustrates the importance of the length of the defect a . It represents the measured tensions for 100 defective tubes with various defect depths. The depths are normalized so that 1 represents the thickness of a tube coating. Using the parametrization presented at the end of the introduction, the 100 points form a space-filling design set for the 9-dimensional cube $(0, 1)^9$.

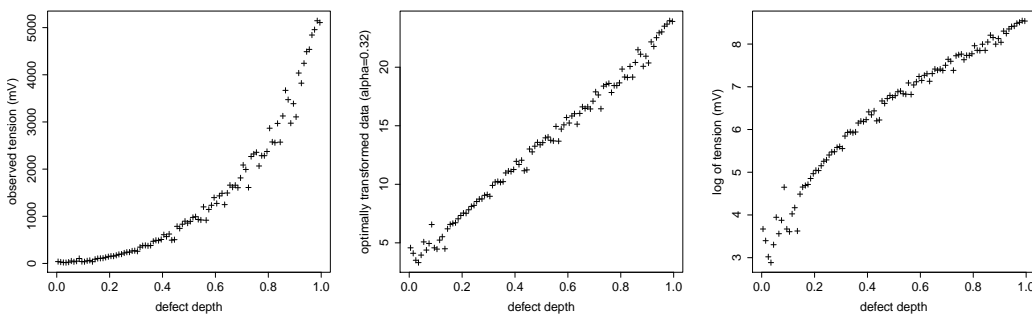


Figure 7.2 – *Left: Measured tension as a mapping of defect depth for 100 defective tubes. The other “nuisance” geometrical parameters are representative of possible geometries. Right: Same data after applying the logarithm function (B_α or C_α with $\alpha = 0$) to all measured tensions. Center: Same data after applying the optimal transformation (C_α with $\alpha = 0.32$) to all measured tensions.*

Because the measured tension is necessarily nonnegative, the Gaussian assumption of the Kriging surrogate model is inadequate. Moreover, the spread becomes greater and greater as the defect depth increases, which contradicts the assumption of stationarity. Because Figure 7.2 (left) gives the impression of being based on the graph of some exponential mapping, our first instinct was to apply the logarithm to the observations. The result is given by Figure 7.2 (right).

Although this transformation seems fruitful, it is clearly too strong for our purpose. While the original data yielded what looked like a strongly convex mapping of defect depth, the log-transformed data yield a somewhat concave mapping of defect depth. Moreover, the spread which was originally increasing with depth seems now to be decreasing with depth. Some intermediate transformation between the identity mapping and the logarithm mapping was needed. A possible choice could have been the Box-Cox power transform family:

$$B_\alpha(t) = \begin{cases} \frac{t^\alpha - 1}{\alpha} & \text{if } \alpha > 0 \\ \log(t) & \text{if } \alpha = 0 \end{cases} \tag{7.3}$$

The Box-Cox power transformation fits our requirements since all mappings B_α for $\alpha \in (0, 1)$ are intermediate transformations between the logarithm ($\alpha = 0$) and the identity mapping ($\alpha = 1$, although the data are uniformly decremented by 1, which is of no consequence).

However, for any $\alpha > 0$, B_α is a bijection from $(0, +\infty)$ to $(-1/\alpha, +\infty)$, whereas we would like a bijection from $(0, +\infty)$ to $(-\infty, +\infty)$. How can the Gaussian assumption be credible otherwise? In the Box-Cox family, only the logarithm mapping ($\alpha = 0$) fits this requirement.

The following alteration to the Box-Cox family is suitable:

$$C_\alpha(t) = \begin{cases} \frac{1}{\alpha} \sinh(\alpha \log(t)) & \text{if } \alpha > 0 \\ \log(t) & \text{if } \alpha = 0 \end{cases} \quad (7.4)$$

Every mapping C_α is a bijection from $(0, +\infty)$ to $(-\infty, +\infty)$. Moreover, C_1 is equivalent to the linear mapping $t \mapsto 0.5t$ when $t \rightarrow \infty$, so this family too can be considered to contain intermediate mappings between the logarithm ($\alpha = 0$) and the identity mapping ($\alpha = 1$).

Figure 7.3 illustrates both Box-Cox and alternative transformation families.

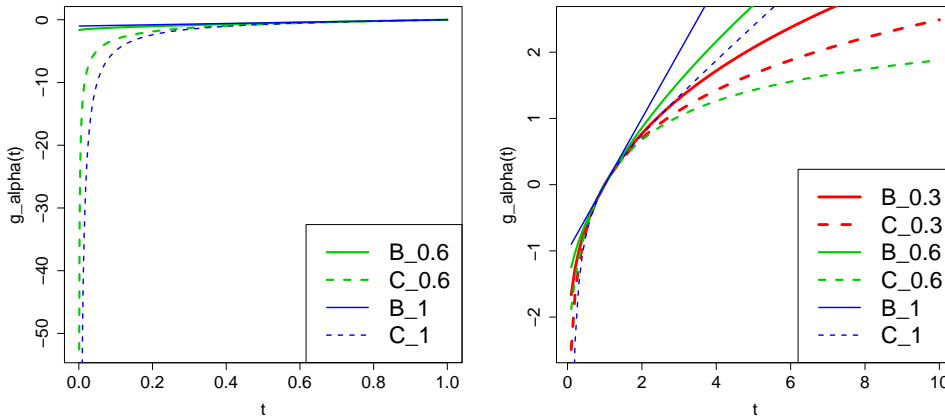


Figure 7.3 – Box-Cox and alternative transformation families for small values of α (left) and greater values of α (right).

We use this transformation family (i.e. $g_\alpha = C_\alpha$ for all $\alpha \in [0, +\infty)$) to apply the Bayesian framework for Trans-Gaussian Kriging described in the previous section. Our mean function space \mathcal{F}_2 is the space of affine functions of the parameter of interest a . This choice reflects the fact that a is the main influence on the output of the computer code. Our basis functions are the constant function of value 1 and the coordinate function $(a, \mathbf{x}) \mapsto a$. Our correlation kernel is the Matérn anisotropic geometric kernel of smoothness $\nu = 5/2$. The transformation family is C_α , $\alpha \in [0, +\infty)$.

Trans-Gaussian Kriging Likelihood

In Trans-Gaussian Kriging, Z is a random field whose distribution is that of \mathcal{Z} conditioned by $\mathcal{Z}(a^{(i)}, \mathbf{x}^{(i)}) = z(a^{(i)}, \mathbf{x}^{(i)})$ for every $i \in \llbracket 1, n \rrbracket$, and \mathbf{z} is the vector $(z(a^{(1)}, \mathbf{x}^{(1)}), \dots, z(a^{(n)}, \mathbf{x}^{(n)}))^\top$.

Denoting by $g_\alpha(\mathbf{z})$ the vector $(g_\alpha(z_1), \dots, g_\alpha(z_n))^\top$, the likelihood of hyperparameters β , σ^2 , θ and α is:

$$L(\mathbf{z} \mid \beta, \sigma^2, \theta, \alpha) = \left(\frac{1}{2\pi\sigma^2}\right)^{\frac{n}{2}} |\Sigma_\theta|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2\sigma^2} (g_\alpha(\mathbf{z}) - \mathbf{H}\beta)^\top \Sigma_\theta^{-1} (g_\alpha(\mathbf{z}) - \mathbf{H}\beta)\right\} \prod_{i=1}^n g'_\alpha(z_i). \quad (7.5)$$

With Trans-Gaussian Kriging, the ‘‘Gaussian’’ parameters β , σ^2 and θ can no longer be interpreted as pertaining to mean function, variance and correlation respectively. They can only be interpreted in relation to the transformation parameter α .

[Box and Cox, 1964] propose estimating the ‘‘Gaussian’’ parameters conditionally to the transformation parameter, and to treat the value of the transformation parameter with maximum likelihood as the true value. This point of view was criticized by Bickel and Doksum [1981], who showed that when using the Box-Cox transformation family, the estimators for the transformation parameter and for β are highly correlated, and they argued for taking this effect into account when performing inference on the parameters. Unfortunately, when the computational budget is too low to be able to make many observations, the Kriging likelihood can be rather flat even without adding a transformation parameter into the mix [Li and Sudjianto, 2005].

These difficulties are the reason why we propose in this chapter using an Objective Bayesian framework to integrate β , σ^2 and θ out of the model – conditionally to α , in keeping with Box and Cox [1964]’s insight. This way, we only need to estimate the parameter α , with the likely consequence that the ‘‘integrated’’ one-dimensional likelihood is not as flat as the regular multi-dimensional likelihood.

Objective Bayesian treatment of the Gaussian parameters

In this subsection, we operate again under the framework of Subsection 7.3. This framework is the standard Kriging setting, which is parametrized by β , σ^2 and θ . The likelihood is given by Equation (7.2).

Following Berger et al. [2001], we make use of the reference prior on the parameters (β, σ^2) conditional to θ : $\pi_R(\beta, \sigma^2 \mid \theta) \propto 1/\sigma^2$. The integrated likelihood $L^1(\mathbf{y} \mid \theta)$ is given by successively integrating β and σ^2 out. To do this, we introduce \mathbf{W} as an $n \times (n - p)$ matrix such that $\mathbf{W}^\top \mathbf{H}$ is the null $(n - p) \times p$ matrix and $\mathbf{W}^\top \mathbf{W}$ is the $(n - p) \times (n - p)$ identity matrix. \mathbf{W} can be computed by performing Standard Value Decomposition (SVD) on \mathbf{H} . The singular vectors corresponding to the null singular values form the columns of \mathbf{W} . We find

$$L^1(\mathbf{y} \mid \theta) = \int L(\mathbf{y} \mid \beta, \sigma^2, \theta) / \sigma^2 d\beta d\sigma^2 = \left(\frac{2\pi^{\frac{n-p}{2}}}{\Gamma(\frac{n-p}{2})}\right)^{-1} |\mathbf{W}^\top \Sigma_\theta \mathbf{W}|^{-\frac{1}{2}} \left(\mathbf{y}^\top \mathbf{W} \left(\mathbf{W}^\top \Sigma_\theta \mathbf{W}\right)^{-1} \mathbf{W}^\top \mathbf{y}\right)^{-\frac{n-p}{2}}. \quad (7.6)$$

We briefly explain how prediction can be performed in our Bayesian Kriging model.

First, assume that all hyperparameters – β , σ^2 and θ – are known. We wish to predict the value taken by the Gaussian process Y at an unobserved point (a_0, \mathbf{x}_0) .

Let \mathbf{P} be the $n \times p$ matrix whose columns are obtained by applying the Gram-Schmidt process to the columns of \mathbf{H} .

Denote by $\mathbf{H}_{0,0}$ the transpose of the vector of length p containing the values of the p basis functions at (a_0, \mathbf{x}_0) , by $\Sigma_{\theta,0,\cdot}$ the $1 \times n$ correlation matrix between the Gaussian process at (a_0, \mathbf{x}_0) and the Gaussian process at the design set, and finally by $\Sigma_{\theta,\cdot,0}$ its transpose.

Define the $1 \times p$ matrix $\mathbf{H}_{0,\cdot} = \mathbf{H}_{0,0} (\mathbf{P}^\top \mathbf{H})^{-1}$ and its transpose $\mathbf{H}_{\cdot,0}$. The following matrix definitions are necessary to express the predictive distribution:

$$\begin{aligned} E_0 &:= \mathbf{H}_{0,\cdot} \mathbf{P}^\top \mathbf{y} \\ S_{\theta,0,0} &:= 1 + \mathbf{H}_{0,\cdot} \mathbf{P}^\top \Sigma_\theta \mathbf{P} \mathbf{H}_{\cdot,0} - \mathbf{H}_{0,\cdot} \mathbf{P}^\top \Sigma_{\theta,\cdot,0} - \Sigma_{\theta,0,\cdot} \mathbf{P} \mathbf{H}_{\cdot,0} \\ S_{\theta,0,\cdot} &:= \left(\mathbf{H}_{0,\cdot} \mathbf{P}^\top \Sigma_\theta - \Sigma_{\theta,0,\cdot} \right) (\mathbf{P} \mathbf{W}) \\ S_{\theta,\cdot,0} &:= S_{\theta,0,\cdot}^\top \end{aligned}$$

We recall here Corollary 6.3:

Proposition 7.1. *The predictive distribution of $Y(a_0, \mathbf{x}_0)$ averaged over both β and σ^2 is the non-standardized Student's t -distribution with $n - p$ degrees of freedom, location parameter $E_0 - S_{\theta,0,\cdot} \mathbf{W} (\mathbf{W}^\top \Sigma_\theta \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{y}$ and scale parameter*

$$\sqrt{\frac{\mathbf{y}^\top \mathbf{W} (\mathbf{W}^\top \Sigma_\theta \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{y}}{n - p} \left\{ S_{\theta,0,0} - S_{\theta,0,\cdot} \mathbf{W} (\mathbf{W}^\top \Sigma_\theta \mathbf{W})^{-1} \mathbf{W}^\top S_{\theta,\cdot,0} \right\}}.$$

Practically speaking, if $n - p$ is large the Student t -distribution can be approximated by a Normal distribution with mean equal to the location parameter and standard deviation equal to the scale parameter.

Further integrating the predictive distribution requires averaging over the Gibbs reference posterior distribution on θ . The Gibbs reference posterior distribution is derived through Objective Bayesian techniques and has nice theoretical as well as practical properties, like invariance by reparametrization and good frequentist performance for prediction intervals, as was seen in Sections 5.4 and 6.5. This integration can be done numerically, by using a sample from the Gibbs reference posterior that can be obtained easily from a Gibbs sampler.

The Likelihood problem

Let us go back to Trans-Gaussian Kriging. In this framework, it is desirable to integrate θ out of the model in order to derive an integrated likelihood of the transformation parameter α only.

The problem with the Gibbs reference posterior approach is that, given there is no actual prior distribution yielding the Gibbs reference posterior, it is not possible to integrate θ out

of the model with likelihood $L^1(\mathbf{y}|\boldsymbol{\theta})$. This is a major hurdle since the actual model (taking into account that \mathbf{y} is the result of a transformation) is:

$$L^1(\mathbf{z}|\boldsymbol{\theta}, \alpha) = \left(\frac{2\pi^{\frac{n-p}{2}}}{\Gamma\left(\frac{n-p}{2}\right)} \right)^{-1} |\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W}|^{-\frac{1}{2}} \left(g_\alpha(\mathbf{z})^\top \mathbf{W} \left(\mathbf{W}^\top \boldsymbol{\Sigma}_\theta \mathbf{W} \right)^{-1} \mathbf{W}^\top g_\alpha(\mathbf{z}) \right)^{-\frac{n-p}{2}} \prod_{i=1}^n g'_\alpha(z_i). \tag{7.7}$$

Because we cannot integrate $\boldsymbol{\theta}$ out, we need to compute some other aggregate of all possible $L^1(\mathbf{z}|\boldsymbol{\theta}, \alpha)$ when $\boldsymbol{\theta}$ varies. The most obvious solution is to average $L^1(\mathbf{z}|\boldsymbol{\theta}, \alpha)$ over the Gibbs reference posterior distribution $\pi_G(\boldsymbol{\theta}|\mathbf{z}, \alpha) := \pi_G(\boldsymbol{\theta}|g_\alpha(\mathbf{z}))$, but this solution gives too much weight to the likelihood: if a prior $\pi_G(\boldsymbol{\theta}|\alpha)$ existed, then this would be equivalent to computing $\int L^1(\mathbf{z}|\boldsymbol{\theta}, \alpha)^2 \pi_G(\boldsymbol{\theta}|\alpha) d\boldsymbol{\theta} / \int L^1(\mathbf{z}|\boldsymbol{\theta}, \alpha) \pi_G(\boldsymbol{\theta}|\alpha) d\boldsymbol{\theta}$.

A second possibility is to compute the Maximum A Posteriori estimator $\widehat{\boldsymbol{\theta}}_{MAP}$ for $\boldsymbol{\theta}$ and then proceed with $L^{MAP}(\mathbf{z}|\alpha) := L^1(\mathbf{z}|\boldsymbol{\theta} = \widehat{\boldsymbol{\theta}}_{MAP}, \alpha)$. The idea is that $\widehat{\boldsymbol{\theta}}_{MAP}$ should be a “typical” value of $\boldsymbol{\theta}$, but this approach unfortunately makes the aggregate dependent on the parametrization chosen for the correlation lengths, which is a bad property for a quantity treated as a likelihood.

A third possibility is

$$L^{LOG}(\mathbf{z}|\alpha) := \exp \left[\int \log\{L^1(\mathbf{z}|\boldsymbol{\theta}, \alpha)\} \pi_G(\boldsymbol{\theta}|\mathbf{z}, \alpha) d\boldsymbol{\theta} \right] \tag{7.8}$$

Compared to L^{MAP} , L^{LOG} has the advantage that it does not depend on the parametrization, insofar as $\boldsymbol{\theta}$ could be replaced by any vector $(h_1(\theta_1), \dots, h_r(\theta_r))^\top$ as long as all h_i ($i \in \llbracket 1, r \rrbracket$) are bijections. Moreover, it does not rely on a particular estimate of the “true” $\boldsymbol{\theta}$, which makes L^{LOG} more robust than L^{MAP} .

Fundamentally though, L^{MAP} and L^{LOG} are justified by a simple heuristic. This heuristic is based on two approximations.

1. Asymptotic: the Gibbs reference posterior distribution is approximated by the probability distribution $\mathcal{N}(\widehat{\boldsymbol{\theta}}_{MAP}; \mathcal{I}(\widehat{\boldsymbol{\theta}}_{MAP})^{-1})$ where $\mathcal{I}(\boldsymbol{\theta})$ is the Fisher information matrix. This means that we assume that the conclusion of the Bernstein - von Misès theorem applies as though we were in an asymptotic framework.
2. Jeffreys: the Gibbs reference posterior is the posterior distribution of $\boldsymbol{\theta}$ corresponding to the Jeffreys-rule prior on $\boldsymbol{\theta}$ denoted by $\pi(\boldsymbol{\theta}|\alpha)$. In order to have simple notations, we do not normalize it – it is possibly not proper anyway – so, denoting by $|\cdot|$ the matrix determinant, $\pi(\boldsymbol{\theta}|\alpha) = |\mathcal{I}(\boldsymbol{\theta})|^{1/2}$.

Define $K_{\mathbf{z}}(\alpha) \in \mathbb{R}$ such that $K_{\mathbf{z}}(\alpha) \int L^1(\mathbf{z}|\boldsymbol{\theta}, \alpha) \pi(\boldsymbol{\theta}|\alpha) d\boldsymbol{\theta} = \prod_{i=1}^n g'_\alpha(z_i)$. This is equivalent to asserting that $\pi_G(\boldsymbol{\theta}|\mathbf{z}, \alpha) \prod_{i=1}^n g'_\alpha(z_i) = K_{\mathbf{z}}(\alpha) L^1(\mathbf{z}|\boldsymbol{\theta}, \alpha) \pi(\boldsymbol{\theta}|\alpha)$. In the following, we show that the two approximations imply that $\prod_{i=1}^n g'_\alpha(z_i) / K_{\mathbf{z}}(\alpha)$ is related to both $L^{MAP}(\mathbf{z}|\alpha)$ and $L^{LOG}(\mathbf{z}|\alpha)$.

First, we have

$$\begin{aligned} & K_{\mathbf{z}}(\alpha)L^1(\mathbf{z}|\boldsymbol{\theta}, \alpha)\pi(\boldsymbol{\theta}|\alpha) \\ & = (2\pi)^{-r/2}|\mathcal{I}(\boldsymbol{\theta})|^{1/2} \exp\left\{-\frac{(\boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}_{MAP})^\top \mathcal{I}(\widehat{\boldsymbol{\theta}}_{MAP})(\boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}_{MAP})}{2}\right\} \prod_{i=1}^n g'_\alpha(z_i). \end{aligned} \quad (7.9)$$

In particular,

$$K_{\mathbf{z}}(\alpha)L^1(\mathbf{z}|\widehat{\boldsymbol{\theta}}_{MAP}, \alpha) = (2\pi)^{-r/2} \prod_{i=1}^n g'_\alpha(z_i). \quad (7.10)$$

So the integrated likelihood $\int L^1(\mathbf{z}|\boldsymbol{\theta}, \alpha)\pi(\boldsymbol{\theta}|\alpha)d\boldsymbol{\theta}$ is proportional to $L^{MAP}(\mathbf{z}|\alpha)$:

$$\int L^1(\mathbf{z}|\boldsymbol{\theta}, \alpha)\pi(\boldsymbol{\theta}|\alpha)d\boldsymbol{\theta} = (2\pi)^{r/2}L^{MAP}(\mathbf{z}|\alpha). \quad (7.11)$$

Furthermore,

$$\int \log L^1(\mathbf{z}|\boldsymbol{\theta}, \alpha)\pi_G(\boldsymbol{\theta}|\mathbf{z}, \alpha)d\boldsymbol{\theta} = \int \log\{K_{\mathbf{z}}(\alpha)L^1(\mathbf{z}|\boldsymbol{\theta}, \alpha)\}\pi_G(\boldsymbol{\theta}|\mathbf{z}, \alpha)d\boldsymbol{\theta} - \log K_{\mathbf{z}}(\alpha) \quad (7.12)$$

Notice that

$$\log\{K_{\mathbf{z}}(\alpha)L^1(\mathbf{z}|\boldsymbol{\theta}, \alpha)\} = -\frac{r}{2}\log(2\pi) - \frac{(\boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}_{MAP})^\top \mathcal{I}(\widehat{\boldsymbol{\theta}}_{MAP})(\boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}_{MAP})}{2} + \sum_{i=1}^n \log g'_\alpha(z_i), \quad (7.13)$$

so

$$\int \log L^1(\mathbf{z}|\boldsymbol{\theta}, \alpha)\pi_G(\boldsymbol{\theta}|\mathbf{z}, \alpha)d\boldsymbol{\theta} = -\frac{r}{2}\log(2\pi) - \frac{r}{2} + \sum_{i=1}^n \log g'_\alpha(z_i) - \log K_{\mathbf{z}}(\alpha). \quad (7.14)$$

Finally, we see that the integrated likelihood $\int L^1(\mathbf{z}|\boldsymbol{\theta}, \alpha)\pi(\boldsymbol{\theta}|\alpha)d\boldsymbol{\theta}$ is also proportional to $L^{LOG}(\mathbf{z}|\alpha)$:

$$\log\left\{\int L^1(\mathbf{z}|\boldsymbol{\theta}, \alpha)\pi(\boldsymbol{\theta}|\alpha)d\boldsymbol{\theta}\right\} = \log L^{LOG}(\mathbf{z}|\alpha) + \frac{r}{2}\log(2\pi) + \frac{r}{2}. \quad (7.15)$$

Interestingly, this heuristic also predicts that

$$\log\{L^{MAP}(\mathbf{z}|\alpha)\} - \log\{L^{LOG}(\mathbf{z}|\alpha)\} = r/2. \quad (7.16)$$

This prediction provides a sanity check for the heuristic, which we use in the application below.

7.4 Industrial Application

Integrating out Kriging hyperparameters

Figure 7.2 shows that α can only reasonably belong to $[0, 1]$, so we endow $[0, 1]$ with the fine grid $0.01 * \llbracket 0, 100 \rrbracket$. For every element α in this grid, we sample 100 points $\boldsymbol{\theta}_\alpha^{(j)} \in (0, +\infty)^r$ from the Gibbs reference posterior distribution. This is done by performing 9000 iterations of the Gibbs algorithm and keeping only one out of 90.

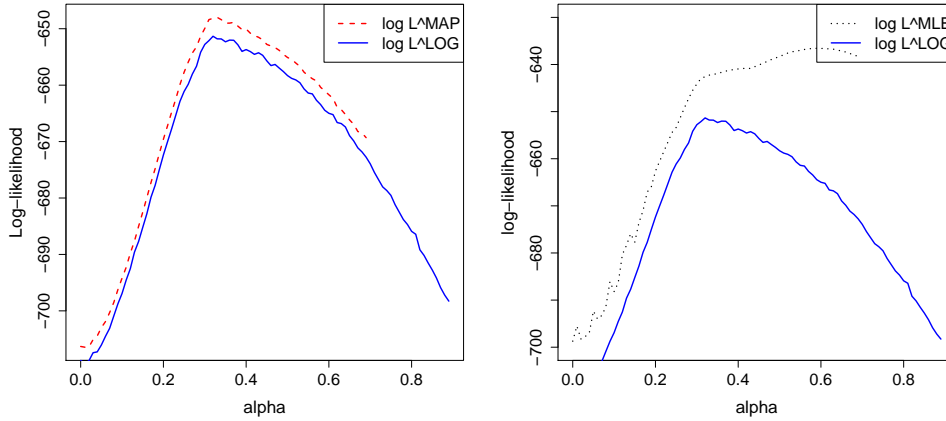


Figure 7.4 – Left: Logarithm of $L^{MAP}(\mathbf{z}|\alpha)$ (red dotted line) and $L^{LOG}(\mathbf{z}|\alpha)$ (blue solid line). L^{MAP} is not represented for greater values of α because the MAP cannot be reliably computed. Right: Logarithm of $L^{LOG}(\mathbf{z}|\alpha)$ (blue solid line) and log-likelihood $L^{MLE}(\mathbf{z}|\alpha)$ if the MLE estimator for θ is used (black dotted line). L^{MAP} and L^{MLE} are not represented for greater values of α , because the MAP cannot be reliably computed and the MLE favors correlation lengths so high that correlation matrices are ill-conditioned.

Figure 7.4 (left) gives the logarithm of the pseudo-likelihoods L^{MAP} and L^{LOG} .

L^{MAP} and L^{LOG} reach their maxima at $\alpha = 0.34$ and $\alpha = 0.32$ respectively. Interestingly, the average value of $L^{MAP}(\mathbf{z}|\alpha) - L^{LOG}(\mathbf{z}|\alpha)$ for $\alpha \in [0.25, 0.55]$ is 3.4 and its standard deviation 0.3. This is quite close to the prediction (7.16) that this difference should be constant, especially considering that the values of both functions spread over an interval of length greater than 50. However, the value 3.4 is substantially lower than the predicted $r/2 = 4.5$, which suggests that not all 9 parameters are influent.

On account of L^{LOG} not depending on the estimate of the MAP, we accept the value of α maximizing L^{LOG} , $\alpha = 0.32$ as the “true” alpha. In this respect, we follow the suggestion of Box and Cox [1964]. Naturally, a completely Bayesian treatment would require defining a prior distribution on α in order to integrate it out. In this case however, both versions of the likelihood are so pronounced that any reasonable prior would have little effect on the posterior distribution. In the interest of completeness though, we do offer in the next section a comparison between the prediction of the Maximum Likelihood approach setting $\alpha = 0.32$ and a “full-Bayesian” approach using the Uniform prior distribution on $[0, 1]$.

The question of whether or not to use a Bayesian approach to deal with α instead of the MLE may be debatable, but it is not when dealing with θ : the Bayesian approach is clearly superior. Figure 7.4 (right) shows the log-likelihood of α when taking θ to be equal to its MLE (dotted curve). It favors high values of α , which means the data are weakly transformed so the MLE on θ favors very high correlation lengths and correlation matrices become ill-conditioned.

Let us consider again the observation data presented in Figure 7.2 (left). If we apply the transformation C_α with $\alpha = 0.32$, we obtain Figure 7.2 (center). Notice that the observations seem to be placed along a straight line. This is no coincidence: it reflects our choice of mean function space \mathcal{F}_2 . The optimal transformation parameter is in first approximation the one that makes the data match our assumption about the mean function.

To confirm this, consider a surrogate model of transformed linear regression. It can be seen as a Trans-Gaussian Kriging model with null correlation lengths. We provide the log-likelihood of the transformation parameter for such a model in Figure 7.4 (right). It reaches its maximum at $\alpha = 0.30$, very near $\alpha = 0.32$ where L^{LOG} reaches its maximum. This shows that the correlation structure parametrized by $\boldsymbol{\theta}$ has little impact on the maximum likelihood of α . $\boldsymbol{\beta}$ and σ^2 having been integrated out of the model, the choice of mean function space \mathcal{F}_p is necessarily the primary explanation for the the likelihood function of α reaching its maximum where it does.

Probability Of defect Detection

With this machinery, we can now return to the question of POD curves. For any point in the $r = 9$ -dimensional input space of C3D, the Trans-Gaussian Kriging surrogate model can provide a probability that the detection threshold $s = 200mV$ is crossed. We need a few definitions.

Definition 7.2. *Let $a \in [0, 1]$. The actual POD, denoted by $POD(a)$, is the probability for a defect of depth a to be detected: $POD(a) = \mathbb{P}(z(a, \mathbf{X}) > s)$*

This probability refers to actual randomness, in the sense that the geometrical characteristics of a defect are considered random. It is a probability in the frequentist sense but cannot be accessed without prohibitive computational costs because it would involve running C3D over a large set of geometries \mathbf{x} .

Definition 7.3. *The surrogate model safety denoted by $SAFE(a, \mathbf{x})$ is the probability, according to the surrogate model, of a particular defect characterized by (a, \mathbf{x}) being detected: $SAFE(a, \mathbf{x}) = \mathbb{P}(Z(a, \mathbf{x}) > s)$.*

Contrarily to actual POD, surrogate model safety does not refer to any randomness in the frequentist meaning but instead expresses the “opinion” of the surrogate model. It represents epistemic uncertainty and is a probability in the Bayesian sense.

Definition 7.4. *Let $a \in [0, 1]$. The mean POD, denoted by $POD_{mean}(a)$, is the average of surrogate model safety over all defects of depth a : $POD_{mean}(a) = \mathbb{E}(SAFE(a, \mathbf{X})) = \mathbb{P}(Z(a, \mathbf{X}) > s)$.*

The mean POD is perhaps the best approximation of the actual POD available to us. However it is difficult to interpret since it aggregates two very different kinds of uncertainty: uncertainty about the defect geometry, which is random in kind, and epistemic uncertainty, which refers to imprecision on the part of the surrogate model.

Definition 7.5. *Let $a \in [0, 1]$ and $\gamma \in [0, 1]$. The POD at safety level γ , denoted by $POD_\gamma(a)$, is the probability for surrogate model safety to be greater or equal to γ : $POD_\gamma(a) = \mathbb{P}(SAFE(a, \mathbf{X}) \geq \gamma)$.*

The POD at safety level γ for a given depth length a is the probability of the surrogate model being confident about the defect being detected, with γ denoting the required confidence level. Its aim is to constrain epistemic uncertainty in order to provide a figure that reflects actual

randomness. Its interpretation is therefore clearer than that of the mean POD. However, it is of interest only if γ is high. In the following, we compute it with $\gamma = 95\%$ and $\gamma = 99\%$.

Computationally speaking, because the defect depth a belongs to the interval $[0, 1]$, we endow this interval with the fine grid $0.01 * \llbracket 0, 100 \rrbracket$. For every value of a in this grid, we generate a 1000-points sample from the probability distribution of \mathbf{X} . For every a therefore, we gather 1000 probabilities of defect detection. With this, we may:

1. compute the mean;
2. count how many are greater or equal to safety levels (95% and 99%).

The first quantity is a Monte-Carlo approximation of the mean POD, the second of the POD at safety level 95% or 99%.

The mean POD and POD at safety level 95% and 99% are drawn in Figure 7.5. On the left, α is taken to be 0.32, the value for which L^{LOG} reaches its maximum. On the right, they are also depicted and accompanied by the curves obtained with integrated α .

Integrated α here means that the Uniform prior distribution on $[0, 1]$ was used for transformation parameter α . Because the likelihood L^{LOG} has a very sharp peak, integration was performed only over the interval $[0.30, 0.37]$. The complement set was deemed to have too low posterior probability to be worth taking into account. Integration was performed by using the rectangle method over $[0.30, 0.37]$ with gap 0.01. This is admittedly a rough approximation, and is intended to serve as a glimpse into a Bayesian treatment of every single parameter of the problem. Figure 7.5 (right) shows that the full-Bayesian method differs only slightly from the Maximum Likelihood method. Interestingly, the full-Bayesian method seems slightly less conservative than the Maximum Likelihood method. We therefore adopt in this case the Maximum Likelihood method for simplicity and conservativeness both.

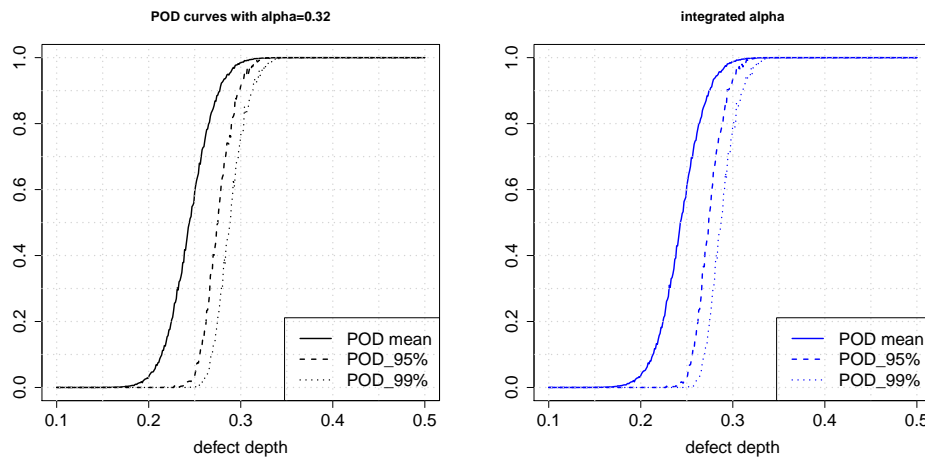


Figure 7.5 – Left: POD curves with $\alpha = 0.32$. Right: POD curves with $\alpha = 0.32$ and for integrated α (points). Since curves with $\alpha = 0.32$ and integrated α are almost confused, it was not necessary to use different graphical styles for the different POD curves derived with integrated α .

7.5 Conclusion

In this chapter, we demonstrated how an Objective Bayesian framework could be applied to a Trans-Gaussian Kriging surrogate model. Contrary to Maximum Likelihood approaches, it makes the likelihood function of the transformation parameter clearly discriminate between all possible candidates. Moreover, it provides a way to naturally incorporate hyperparameter uncertainty into the prediction (at least as far as “Gaussian parameters” are concerned, because one of our findings was precisely that there is not much uncertainty about the transformation parameter).

This is especially useful in the context of calibration of NDT techniques by numerical simulation. It makes a clear distinction between randomness and epistemic uncertainty possible. This distinction can then be incorporated into Probability Of defect Detection (POD) curves.

Concerning the particular problem of detecting defects in Steam Generator tubes, our framework provides a theoretically sound solution for estimating surrogate model uncertainty. This makes interpretation of results easier and thereby increases the reliability of the results of the study.

Conclusion

What was done

This thesis is the result of the combination of methods originating from geostatistics, Objective Bayesian statistics and from Markov chain theory. All three domains are part of the solution to an industrial problem put forward by EDF.

This problem was how to compute Probability Of defect Detection (POD) curves for non-destructive testing of defects in steam generator tubes. Because a computer simulation of the procedure was available, such curves could theoretically have been obtained by repeated calls to the code under a Monte-Carlo procedure. Budgetary constraints made such an approach untractable, however.

It is precisely to deal with such cases that surrogate models like Kriging exist. However, the scarcity of available data made reliable estimation of Kriging hyperparameters impossible. Hence the use of Objective Bayesian techniques to circumvent the need to estimate these parameters in the first place.

The first part of the present dissertation describes the tools that were the basis of the solution: Kriging and reference analysis. It recalls how Berger et al. [2001] combined these tools to provide a complete Objective Bayesian analysis of the Kriging model with isotropic covariance kernel and fixes problems in the original proof.

The third part of the dissertation builds on Berger et al. [2001]’s success to extend the analysis to Kriging models with anisotropic kernels. It shows how such a Bayesian analysis makes the derivation of POD curves natural even in the absence of reliable estimate for the Kriging hyperparameters. It also demonstrates how this approach is compatible with the elementary non-Gaussian technique called “trans-Gaussian Kriging” which is basically regular Kriging with an added parametric transformation. Even though complexity increases with each new parameter, it sketches a full-Bayesian analysis of the full model, including the transformation parameter.

This would not have been possible without the notion of optimal compromise between potentially incompatible conditional distributions, which is introduced in the second part. This notion is of great practical interest because it allows Gibbs sampling and thereby makes the full-Bayesian solution tractable. An alternative was developed recently [Gu et al., 2018]: it uses the reference prior obtained by grouping all covariance parameters in the reference prior algorithm. Practical implementations however rely on the Maximum A Posteriori estimator instead of conducting a full-Bayesian analysis.

What remains to be done

There is much that remains to be understood about the optimal compromise between incompatible conditional distributions. By definition, it fits an optimality criterion, but its properties remain largely unknown. What is the precise link between each of the incompatible conditionals and the corresponding conditional in the compromise? My intuition is that the optimal compromise tends to blur dependence relationships, but this idea would need to be quantified.

Is it possible to quantify how much each conditional is taken into account in the compromise? Could changes to the definition allow for the regulation of the importance of each conditional in the definition of the compromise?

For example, an intuitive idea would be to modify the scanning probabilities. If we wanted one of the conditionals to be better taken into account, we could alter the Gibbs algorithm and make the probability of choosing it greater.

But whether the impact of a conditional predominantly depends on this probability or on other factors is unclear. A rarely chosen conditional could still have a tremendous impact on one of the marginals of the resulting joint distribution and thereby impact the joint distribution itself.

And above all, a compromise being optimal does not actually mean that it is *good*. Criteria should be developed to quantify how satisfying a compromise is. In the application presented in the dissertation, results in terms of frequentist behavior were surprisingly good, but they remain empirical rather than formal.

Bibliography

- M. Abramowitz and I. A. Stegun, editors. *Handbook of mathematical functions with formulas, graphs, and mathematical tables*, volume 55 of *Applied Mathematics Series*. National Bureau of Standards, 1964.
- R. J. Adler. *An introduction to continuity, extrema, and related topics for general Gaussian processes*. Hayward, CA: Institute of Mathematical Statistics, 1990.
- E. Anderes. On the consistent separation of scale and variance for Gaussian random fields. *Annals of Statistics*, 38(2):870–893, 2010.
- B. C. Arnold, E. Castillo, and J. M. Sarabia. Conditionally specified distributions: an introduction. *Statistical Science*, 16(3):268–269, 2001.
- F. Bachoc. *Parametric estimation of covariance function in Gaussian-process based Kriging models. Application to uncertainty quantification for computer models*. PhD thesis, Université Paris Diderot, 2013a.
- F. Bachoc. Cross validation and maximum likelihood estimations of hyper-parameters of gaussian processes with model misspecification. *Computational Statistics & Data Analysis*, 66:55–69, 2013b.
- F. Bachoc, G. Bois, J. Garnier, and J.-M. Martinez. Calibration and improved prediction of computer models by Universal Kriging. *Nuclear Science and Engineering*, 176(1):81–97, 2014.
- F. Bachoc, E. Contal, H. Maatouk, and D. Rulli ere. Gaussian processes for computer experiments. *ESAIM: Proceedings and Surveys*, 60:163–179, 2017a.
- F. Bachoc, F. Gamboa, J.-M. Loubes, and N. Venet. A Gaussian Process Regression Model for Distribution Inputs. *IEEE Transactions on Information Theory*, 2017b.
- Richard F Bass. *Probabilistic techniques in analysis*. Springer Science & Business Media, 1995.
- J. Bect, F. Bachoc, and D. Ginsbourger. A supermartingale approach to Gaussian process based sequential design of experiments. *arXiv:1608.01118*, 2016.
- J. Berger. The case for objective bayesian analysis. *Bayesian analysis*, 1(3):385–402, 2006.
- J. O. Berger and J. M. Bernardo. On the Development of Reference Priors. *Bayesian statistics*, 4(4):35–60, 1992.

- J. O. Berger, V. De Oliveira, and B. Sansó. Objective Bayesian analysis of spatially correlated data. *Journal of the American Statistical Association*, 96(456):1361–1374, 2001.
- J. O. Berger, J. M. Bernardo, and D. Sun. The formal definition of reference priors. *Annals of Statistics*, 37(2):905–938, 2009.
- James O Berger, Jose M Bernardo, and Dongchu Sun. Overall objective priors. *Bayesian Analysis*, 10(1):189–221, 2015.
- J. M. Bernardo. Reference Posterior Distributions for Bayesian Inference. *Journal of the Royal Statistical Society. Series B (Methodological)*, 41(2):113–147, 1979a.
- J. M. Bernardo. Expected information as expected utility. *The Annals of Statistics*, 7(3):686–690, 1979b.
- J. M. Bernardo. Reference analysis. In D. Dey and C. Rao, editors, *Handbook of statistics*, volume 25, pages 17–90. Elsevier, 2005.
- J. M. Bernardo. Integrated objective Bayesian estimation and hypothesis testing. In J. M. Bernardo, M. J. Bayarri, and J. O. Berger, editors, *Bayesian statistics 9*, pages 1–68. Oxford University Press, 2011.
- P. J. Bickel and K. A. Doksum. An analysis of transformations revisited. *Journal of the American Statistical Association*, 76(374):296–311, 1981.
- P. Billingsley. *Probability and measure* (Third Edition). Wiley, New York, 1995.
- M. Binois, D. Ginsbourger, and O. Roustant. Quantifying uncertainty on Pareto fronts with Gaussian process conditional simulations. *European journal of operational research*, 243(2):386–394, 2015.
- S. Bochner. *Vorlesungen über Fouriersche Integrale*, volume 12. Akad. Verl.-Ges., 1932.
- G. E. P. Box and D. R. Cox. An analysis of transformations. *Journal of the Royal Statistical Society. Series B (Methodological)*, 26:211–252, 1964.
- R. L. Bras and I. Rodriguez-Iturbe. *Random functions and hydrology*. Courier corporation, 1985.
- C. Chevalier, J. Bect, D. Ginsbourger, E. Vazquez, V. Picheny, and Y. Richet. Fast parallel kriging-based stepwise uncertainty reduction with application to the identification of an excursion set. *Technometrics*, 56(4):455–465, 2014.
- G. Christakos. *Random Field Models in Earth Sciences*. Academic press, New York, 1992.
- B. S. Clarke and A. R. Barron. Jeffreys’ prior is asymptotically least favorable under entropy risk. *Journal of Statistical planning and Inference*, 41(1):37–60, 1994.
- G. Damblin, M. Couplet, and B. Iooss. Numerical studies of space-filling designs: optimization of Latin Hypercube Samples and subprojection properties. *Journal of Simulation*, 7(4):276–289, 2013.

- A. P. Dawid. Invariant prior distributions. In S. Kotz and N. L. Johnson, editors, *Encyclopedia of statistical sciences*, pages 228–236. John Wiley, New York, 1983.
- A. P. Dawid and S. L. Lauritzen. Compatible prior distributions. *Bayesian methods with applications to science, policy and official statistics*, pages 109–118, 2001.
- V. De Oliveira, B. Kedem, and D. A. Short. Bayesian Prediction of Transformed Gaussian Random Fields. *Journal of the American Statistical Association*, 92:1422–1433, 1997.
- P. Druilhet and J.-M. Marin. Invariant HPD and MAP based on Jeffreys measure. *Bayesian Analysis*, 2(4):681–692, 2007.
- R. M. Dudley. *Real Analysis and Probability*. Chapman and Hall/CRC, 1989.
- S. E. Fienberg. Does it make sense to be an "objective bayesian"?(comment on articles by berger and by goldstein). *Bayesian Analysis*, 1(3):429–432, 2006.
- Andrew Gelman and TE Raghunathan. Comment on “Conditionally specified distributions: an introduction” by b.c. arnold, e. castillo and j.m. sarabia. *Statistical Science*, 16(3): 268–269, 2001.
- Andrew Gelman, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin. *Bayesian Data Analysis, Third Edition*. CRC Press, Chapman & Hall, 2013.
- M. Ghosh. Objective Priors: An Introduction for Frequentists. *Statistical Science*, 26(2): 187–202, 2011.
- I. I. Gihman and A. V. Skorohod. *The theory of stochastic processes*, volume 1. Springer-Verlag, Berlin, 1974.
- I. J. Good. A Derivation of the Probabilistic Explication of Information. *Journal of the Royal Statistical Society. Series B (Methodological)*, 28(3):578–581, 1966.
- P. Goovaerts. *Geostatistics for Natural Resources Evaluation*. Oxford University Press, New York, 1997.
- J.C. Gower. Properties of Euclidean and non-Euclidean distance matrices. *Linear Algebra and its Applications*, 67:81–97, 1985.
- M. Gu. *Robust Uncertainty Quantification and Scalable Computation for Computer Models with Massive Output*. PhD thesis, Duke University, 2016.
- M. Gu, X. Wang, and J. O. Berger. Robust gaussian stochastic process emulation. *Annals of Statistics, In Press*, 2018.
- M. S. Handcock and M. L. Stein. A Bayesian Analysis of Kriging. *Technometrics*, 35:403–410, 1993.
- M. S. Handcock and J. R. Wallis. An approach to statistical spatial-temporal modeling of meteorological fields (with discussion). *Journal of the American Statistical Association*, 89 (426):368–390, 1994.

- J. Hartigan. Invariant prior distributions. *The Annals of Mathematical Statistics*, 35(2): 836–845, 1964.
- D. Heckerman, D. M. Chickering, C. Meek, R. Rounthwaite, and C. Kadie. Dependency networks for inference, collaborative filtering, and data visualization. *Journal of Machine Learning Research*, 1:49–75, 2000.
- H. Hotelling. New light on the correlation coefficient and its transforms. *Journal of the Royal Statistical Society. Series B (Methodological)*, 15(2):193–232, 1953.
- B. Iooss and L. Le Gratiet. Uncertainty and sensitivity analysis of functional risk curves based on Gaussian processes. *Reliability Engineering & System Safety*, 2017.
- E. H. Isaaks and M. R. Srivastava. *Applied geostatistics*. Oxford University Press, New York, 1989.
- E. T. Jaynes. Prior probabilities. *IEEE Transactions on systems science and cybernetics*, 4(3):227–241, 1968.
- E. T. Jaynes. On the rationale of maximum-entropy methods. *Proceedings of the IEEE*, 70(9):939–952, 1982.
- H. Jeffreys. *Theory of Probability*. London: Oxford University Press, 3 edition, 1961.
- A. G. Journel and Ch. J. Huijbregts. *Mining geostatistics*. Academic press, New York, 1978.
- K. Kamary and C. P. Robert. Reflecting about selecting noninformative priors. *Journal of Applied and Computational Mathematics*, 3(5):1–7, 2014.
- K. Kamary, K. Mengersen, C. P. Robert, and J. Rousseau. Testing hypotheses via a mixture estimation model. *arXiv:1412.2044*, 2014.
- R. E. Kass and L. Wasserman. The selection of prior distributions by formal rules. *Journal of the American Statistical Association*, 91(435):1343–1370, 1996.
- H. Kazianka and J. Pilz. Objective bayesian analysis of spatial data with uncertain nugget and range parameters. *Canadian Journal of Statistics*, 40(2):304–327, 2012.
- M. C. Kennedy and A. O’Hagan. Bayesian calibration of computer models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 63(3):425–464, 2001.
- D. Khoshnevisan. *Multiparameter Processes: an Introduction to Random Fields*. Springer Science & Business Media, 2002.
- P. K. Kitadinis. *Introduction to Geostatistics: Applications in Hydrogeology*. Cambridge University Press, New York, 1997.
- K. L. Kuo and Y. J. Wang. Pseudo-Gibbs sampler for discrete conditional distributions. *Annals of the Institute of Statistical Mathematics*, pages 1–13, 2017.
- K.-L. Kuo, C.-C. Song, and T. J. Jiang. Exactly and almost compatible joint distributions for high-dimensional discrete conditional distributions. *Journal of Multivariate Analysis*, 157:115–123, 2017.

- P. S. Laplace. *Essai philosophique sur les probabilités*. Vve Courcier, Paris, 1814.
- L. Le Gratiet and J. Garnier. Recursive co-Kriging model for design of computer experiments with multiple levels of fidelity. *International Journal for Uncertainty Quantification*, 4(5), 2014.
- L. Le Gratiet, B. Iooss, G. Blatman, T. Browne, S. Cordeiro, and B. Goursaud. Model Assisted Probability of Detection curves: New statistical tools and progressive methodology. *Journal of Nondestructive Evaluation*, 36(1):8, 2017.
- Loic Le Gratiet. Bayesian analysis of hierarchical multifidelity codes. *SIAM/ASA Journal on Uncertainty Quantification*, 1(1):244–269, 2013.
- P. M. Lee. On the axioms of information theory. *The Annals of Mathematical Statistics*, 35(1):415–418, 1964.
- R. Li and A. Sudjianto. Analysis of computer experiments using penalized likelihood in Gaussian Kriging models. *Technometrics*, 47(2):111–120, 2005.
- D. V. Lindley. On a measure of the information provided by an experiment. *The Annals of Mathematical Statistics*, 27(4):986–1005, 1956.
- D. V. Lindley. Fiducial distributions and bayes' theorem. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 102–107, 1958.
- B. Liseo. The elimination of nuisance parameters. In D. Dey and C. Rao, editors, *Handbook of Statistics*, volume 25, chapter 7, pages 193–219. Elsevier-Sciences, 2006.
- J.-M. Marin, K. Mengersen, and C. P. Robert. Bayesian modelling and inference on mixtures of distributions. In D. Dey and C. Rao, editors, *Handbook of statistics*, volume 25, pages 459–507. Elsevier-Sciences, 2005.
- J.-M. Marin, P. Pudlo, C. P. Robert, and R. J. Ryder. Approximate Bayesian computational methods. *Statistics and Computing*, 22(6):1167–1180, 2012.
- J.-M. Marin, N. S. Pillai, C. P. Robert, and J. Rousseau. Relevant statistics for Bayesian model choice. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(5):833–859, 2014.
- G. Matheron. Krigeage d'un panneau rectangulaire par sa périphérie. *Note géostatistique*, 28, 1960.
- B. Matérn. *Spatial Variation*. Springer-Verlag, Berlin, 2nd edition, 1986.
- L. Maurice, V. Costan, and P. Thomas. Axial probe eddy current inspection of steam generator tubes near anti-vibration bars: performance evaluation using finite element modeling. In *Proceedings of JRC-NDE, Cannes*, pages 638–644, 2013.
- K. Mengersen, C. P. Robert, and M. Titterton. *Mixtures: estimation and applications*, volume 896. John Wiley & Sons, 2011.
- J. Muré. Optimal compromise between incompatible conditional probability distributions and Objective Bayesian Kriging. <https://arxiv.org/pdf/1703.07233>, 2017.

- J. Muré. A Comprehensive Bayesian Treatment of the Universal Kriging parameters with Matérn correlation kernels. <https://arxiv.org/pdf/1801.01007>, 2018.
- J. Muré. Propriety of the reference posterior distribution in Gaussian Process regression. <https://arxiv.org/pdf/1805.08992>, 2018a.
- J. Muré. Trans-Gaussian Kriging in a Bayesian framework : a case study. <https://arxiv.org/pdf/1805.09038>, 2018b.
- L. Nachbin. *The Haar Integral*. New York: van Nostrand, 1965.
- E. Padonou and O. Roustant. Polar Gaussian Processes and Experimental Designs in Circular Domains. *SIAM/ASA Journal on Uncertainty Quantification*, 4(1):1014–1033, 2016.
- R. Paulo. Default priors for Gaussian processes. *Annals of Statistics*, 33(2):556–582, 2005.
- C. Perrin, G. and Soize, S. Marque-Pucheu, and J. Garnier. Nested polynomial trends for the improvement of Gaussian process-based predictors. *Journal of Computational Physics*, 346:389–402, 2017.
- L. Pronzato. Minimax and maximin space-filling designs: some properties and methods for construction. *Journal de la Société Française de Statistique*, 158(1):7–36, 2017.
- L. Pronzato and M.-J. Rendas. Bayesian local Kriging. *Technometrics*, 59(3):293–304, 2017.
- C. E. Rasmussen and C. K. I. Williams. *Gaussian processes for machine learning*. MIT Press, 2006.
- C. Ren, D. Sun, and C. He. Objective bayesian analysis for a spatial model with nugget effects. *Journal of Statistical Planning and Inference*, 142(7):1933–1946, 2012.
- C. Ren, D. Sun, and S. K. Sahu. Objective bayesian analysis of spatial models with separable correlation functions. *Canadian Journal of Statistics*, 41(3):488–507, 2013.
- P. Ressel. De Finetti-type theorems: an analytical approach. *The Annals of Probability*, 13(3):898–922, 1985.
- C. P. Robert. *The Bayesian Choice : From Decision-Theoretic Foundations to Computational Implementation*. Springer-Verlag, New York, 2007.
- C. P. Robert and J. Rousseau. How Principled and Practical Are Penalised Complexity Priors? *Statistical Science*, 32(1):36–40, 2017.
- C. P. Robert, N. Chopin, and J. Rousseau. Harold Jeffreys’s Theory of Probability Revisited. *Statistical Science*, 24(2):141–172, 2009.
- C. P. Robert, J.-M. Marin, and J. Rousseau. Bayesian inference and computation. *Handbook of statistical systems biology*. West Sussex: John Wiley & Sons, Ltd, pages 39–65, 2011.
- O. Roustant, E. Padonou, Y. Deville, A. Clément, G. Perrin, J. Giorla, and H. Wynn. Group kernels for Gaussian process metamodels with categorical inputs. *arXiv:1802.02368*, 2018.

- G. Roverato, A. and Consonni. Compatible prior distributions for directed acyclic graph models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 66(1): 47–61, 2004.
- D. Rullière, N. Durrande, F. Bachoc, and C. Chevalier. Nested Kriging predictions for datasets with a large number of observations. *Statistics and Computing*, 28(4):849–867, 2018.
- T. J. Santner, B. J. Williams, and W. I. Notz. *The Design and Analysis of Computer Experiments*. Springer-Verlag, New York, 2003.
- I.T. Schoenberg. On certain Metric Spaces arising from Euclidean Spaces by a change of metric and their imbedding in Hilbert Space. *Annals of Mathematics*, 38(4):787–793, 1937.
- B. Schölkopf and A. J. Smola. *Learning with Kernels: Support with Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, Cambridge, MA, 2012.
- J. W. Seaman III, J. W. Seaman Jr, and J. D. Stamey. Hidden dangers of specifying noninformative priors. *The American Statistician*, 66(2):77–84, 2012.
- Teddy Seidenfeld. Entropy and uncertainty. *Philosophy of Science*, 53(4):467–491, 1986.
- C. E. Shannon. A Mathematical Theory of Communication. *Bell System Technical Journal*, 27(3):379–423, 1948.
- D. Simpson, H. Rue, A. Riebler, T. G. Martins, and S. H. Sørbye. Penalising Model Component Complexity: A Principled, Practical Approach to Constructing Priors. *Statistical science*, 32(1):1–28, 2017.
- E. Snelson, C. E. Rasmussen, and Z. Ghahramani. Warped gaussian processes. In *Advances in neural information processing systems*, volume 16, pages 337–344, 2004.
- M. L. Stein. *Interpolation of Spatial Data. Some Theory for Kriging*. Springer Series in Statistics. Springer-Verlag, New York, 1999.
- P. Thomas, B. Goursaud, L. Maurice, and S. Cordeiro. Eddy-current non destructive testing with the finite element tool Code_Carmel3d. In *11th International Conference on Non-destructive Evaluation, Jeju*, 2015.
- F. Tuyl, R. Gerlach, and K. Mengersen. Consensus priors for multinomial and binomial ratios. *Journal of Statistical Theory and Practice*, 10(4):736–754, 2016.
- H. Wackernagel. *Multivariate Geostatistics*. Springer-Verlag, Berlin, 1995.
- J. Watson and C. Holmes. Approximate models and robust decisions. *Statistical Science*, 31(4):465–489, 2016.
- B. L. Welch and H. W. Peers. On formulae for confidence points based on integrals of weighted likelihoods. *Journal of the Royal Statistical Society. Series B (Methodological)*, 25(2):318–329, 1963.
- W. Xueou, D. J. Nott, C. C. Drovandi, K. Mengersen, and M. Evans. Using history matching for prior choice. *Technometrics*, 2018.

-
- A. M. Yaglom. *Correlation Theory of Stationary and Related Random Functions 1. Basic Results*. Springer-Verlag, New York, 1987.
- R. Yang and J. O. Berger. *A catalog of noninformative priors*. Institute of Statistics and Decision Sciences, Duke University, 1996.
- B. Yet and W. Marsh. Compatible and incompatible abstractions in Bayesian networks. *Knowledge-Based Systems*, 62:84–97, 2014.
- H. Zhang. Inconsistent Estimation and Asymptotically Equal Interpolations in Model-based Geostatistics. *Journal of the American Statistical Association*, 99(465):250–261, 2004.