



**HAL**  
open science

## Nonparametric estimation of covariance functions by model selection

Jérémie Bigot, Rolando Biscay, Jean-Michel Loubes, Lilian Muniz Alvarez

► **To cite this version:**

Jérémie Bigot, Rolando Biscay, Jean-Michel Loubes, Lilian Muniz Alvarez. Nonparametric estimation of covariance functions by model selection. 2009. hal-00420301

**HAL Id: hal-00420301**

**<https://hal.science/hal-00420301>**

Preprint submitted on 28 Sep 2009

**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.

# Nonparametric estimation of covariance functions by model selection

J. Bigot, R. Biscay, J-M. Loubes and L. Muniz

## Abstract

We propose a model selection approach for covariance estimation of a multi-dimensional stochastic process. Under very general assumptions, observing i.i.d replications of the process at fixed observation points, we construct an estimator of the covariance function by expanding the process onto a collection of basis functions. We study the non asymptotic property of this estimate and give a tractable way of selecting the best estimator among a possible set of candidates. The optimality of the procedure is proved via an oracle inequality which warrants that the best model is selected.

**Keywords:** Model Selection, Covariance Estimation, Stochastic process, Basis expansion, Oracle inequality.

**Subject Class. MSC-2000:** 62G05, 62G20 .

## 1 Introduction

Covariance estimation is a fundamental issue in inference for stochastic processes with many applications, ranging from hydroscience, geostatistics, financial series or epidemiology for instance (we refer to [Ste99], [Jou77] or [Cre93] for general references for applications). Parametric methods have been extensively studied in the statistical literature (see [Cre93] for a review) while nonparametric procedure have received a growing attention along the last decades. One of the main issue in this framework is to impose that the estimator is also a covariance function, preventing the direct use of usual nonparametric statistical methods. In this paper, we propose to use a model selection procedure to construct a nonparametric estimator of the covariance function of a stochastic process under general assumptions for the process. In particular we will not assume Gaussianity nor stationarity.

Consider a stochastic process  $X(t)$  with values in  $\mathbb{R}$ , indexed by  $t \in T$ , a subset of  $\mathbb{R}^d$ ,  $d \in \mathbb{N}$ . Throughout the paper, we assume that  $X$  has finite covariance  $\sigma(s, t) = cov(X(s), X(t)) < +\infty$  for all  $s, t \in T$  and, for sake of simplicity, zero mean  $\mathbb{E}(X(t)) = 0$  for all  $t \in T$ . The observations are  $X_i(t_j)$  for  $i = 1, \dots, N$ ,  $j = 1, \dots, n$ , where the observation points  $t_1, \dots, t_n \in T$  are fixed, and  $X_1, \dots, X_N$  are independent copies of the process  $X$ . Our aim is to build a nonparametric estimator of its covariance.

Functional approximations of the processes  $X_1, \dots, X_N$  from data  $(X_i(t_j))$  are involved in covariance function estimation. When dealing with functional data analysis (see, e.g., [RS05]), smoothing the processes  $X_1, \dots, X_N$  is sometimes carried out as a first step before computing the empirical covariance such as spline interpolation for example (see for instance in [ETA03]) or projection onto a general finite basis. Let  $\mathbf{x}_i = (X_i(t_1), \dots, X_i(t_n))^T$

be the vector of observations at the points  $t_1, \dots, t_n$  with  $i \in \{1, \dots, N\}$ . Let  $\{g_\lambda\}_{\lambda \in \mathcal{M}}$  be a collection of (usually linearly independent but not always) functions  $g_\lambda : T \rightarrow \mathbb{R}$  where  $\mathcal{M}$  denote a generic countable set of indices. Then, let  $(m) \subset \mathcal{M}$  be a subset of indices of size  $m \in \mathbb{N}$  and define the  $n \times m$  matrix  $\mathbf{G}$  with entries  $g_{j\lambda} = g_\lambda(t_j)$ ,  $j = 1, \dots, n$ ,  $\lambda \in (m)$ .  $\mathbf{G}$  will be called the design matrix corresponding to the set of basis functions indexed by  $(m)$ .

In such setting, usual covariance estimation is a two-step procedure: first, for each  $i = 1, \dots, N$ , fit the regression model

$$\mathbf{x}_i = \mathbf{G}\mathbf{a}_i + \epsilon_i \quad (1.1)$$

(by least squares or regularized least squares), where  $\epsilon_i$  are random vectors in  $\mathbb{R}^n$ , to obtain estimates  $\hat{\mathbf{a}}_i = (\hat{a}_{i,\lambda})_{\lambda \in (m)} \in \mathbb{R}^m$  of  $\mathbf{a}_i$  where in the case of standard least squares estimation (assuming for simplicity that  $\mathbf{G}^T \mathbf{G}$  is invertible)

$$\hat{\mathbf{a}}_i = (\mathbf{G}^T \mathbf{G})^{-1} \mathbf{G}^T \mathbf{x}_i, i = 1, \dots, N.$$

Then, estimation of the covariance is given by computing the following estimate

$$\hat{\Sigma} = \mathbf{G} \hat{\Psi} \mathbf{G}^T, \quad (1.2)$$

where

$$\hat{\Psi} = \frac{1}{N} \sum_{i=1}^N \hat{\mathbf{a}}_i \hat{\mathbf{a}}_i^T = (\mathbf{G}^T \mathbf{G})^{-1} \mathbf{G}^T \left( \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{G} (\mathbf{G}^T \mathbf{G})^{-1}. \quad (1.3)$$

This corresponds to approximate the process  $X_i$  by a truncated process  $\tilde{X}_i$  defined as

$$\tilde{X}_i(t) = \sum_{\lambda \in (m)} \hat{a}_{i,\lambda} g_\lambda(t), i = 1, \dots, N,$$

and to choose the empirical covariance of  $\tilde{X}$  as an estimator of the covariance of  $X$ , defined by

$$\hat{\sigma}(s, t) = \frac{1}{N} \sum_{i=1}^N \tilde{X}_i(s) \tilde{X}_i(t).$$

In this paper we propose to view the estimator (1.2) as the covariance obtained by considering a least squares estimator in the following matrix regression model

$$\mathbf{x}_i \mathbf{x}_i^T = \mathbf{G} \Psi \mathbf{G}^T + \mathbf{U}_i, \quad i = 1, \dots, N, \quad (1.4)$$

where  $\Psi$  is a symmetric matrix and  $\mathbf{U}_i$  are i.i.d matrix errors. Fitting the models (1.1) and (1.4) by least squares naturally leads to the definition of different contrast and risk functions as the estimation is not performed in the same space ( $\mathbb{R}^m$  for model (1.1) and  $\mathbb{R}^{m \times m}$  for model (1.4)). By choosing an appropriate loss function, least squares estimation in model (1.4) also leads to the natural estimate (1.2) derived from least square estimation in model (1.1). However, the problem of model selection, i.e. choosing an appropriate data-based subset of indices  $(m) \in \mathcal{M}$ , is very distinct in model (1.1) and model (1.4). Indeed, model selection for (1.1) depends on the variability of the vectors  $\mathbf{x}_i$ 's while for (1.4) it depends on the variability of the matrices  $\mathbf{x}_i \mathbf{x}_i^T$ 's. One of the main contributions of this paper is to show that considering model (1.4) enables to handle a large variety of cases and to build an optimal model selection estimator of the covariance without too

strong assumptions on the model. Moreover it will be shown that considering model (1.4) leads to the estimator  $\hat{\Psi}$  (1.3) which is guaranteed to be in the class of definite non negative matrices and thus to a proper covariance matrix  $\hat{\Sigma} = \mathbf{G}\hat{\Psi}\mathbf{G}^T$ .

A similar method has been developed for smooth interpolation of covariance functions in [BJG95]. However, this paper is restricted to basis functions that are determined by reproducing kernels in suitable Hilbert spaces. Furthermore, a matrix metric different from (though related to) the Frobenius matrix norm is adopted as a fitting criterion. Similar ideas are tackled in [MP08]. These authors deal with the estimation of  $\Sigma$  within the covariance class  $\Gamma = \mathbf{G}\Psi\mathbf{G}^T$  induced by an orthogonal wavelet expansion. However, their fitting criterion is not general since they choose the Gaussian likelihood as a contrast function, and thus their method requires specific distributional assumptions. We also point out that computation of the Gaussian likelihood requires inversion of  $\mathbf{G}\Psi\mathbf{G}^T$ , which is not directly feasible if  $\text{rank}(\mathbf{G}) < n$  or some diagonal entities of the definite non negative (d.n.n) matrix  $\Psi$  are zero.

Hence, to our knowledge, no previous work has proposed to use the matrix regression model (1.4) under general moments assumptions of the process  $X$  using a general basis expansion for nonparametric covariance function estimation.

The paper then falls into the following parts. The description of the statistical framework of the matrix regression is given in Section 2. Section 2 is devoted to the main statistical results. Namely we study the behavior of the estimator for a fixed model in Section 2.1 while Section 2.2 deals with the model selection procedure and provide the oracle inequality. Section 3 states a concentration inequality that is used in all the paper, while the proofs are postponed to a technical Appendix .

## 2 Nonparametric Model selection for Covariance estimation

Recall that  $X = (X(t))_{t \in T}$  is an  $\mathbb{R}$ -valued stochastic process, where  $T$  denotes some subset of  $\mathbb{R}^d$ ,  $d \in \mathbb{N}$ . Assume that  $X$  has finite moments up to order 4, and zero mean, i.e  $\mathbb{E}(X(t)) = 0$  for all  $t \in T$ . The covariance function of  $X$  is denoted by  $\sigma(s, t) = \text{cov}(X(s), X(t))$  for  $s, t \in T$  and recall that  $X_1, \dots, X_N$  are independent copies of the process  $X$ .

In this work, we observe at different observation points  $t_1, \dots, t_n \in T$  these independent copies of the process, denoted by  $X_i(t_j)$ , with  $i = 1, \dots, N$ ,  $j = 1, \dots, n$ . Recall that  $\mathbf{x}_i = (X_i(t_1), \dots, X_i(t_n))^T$  is the vector of observations at the points  $t_1, \dots, t_n$  for each  $i = 1, \dots, N$ . The matrix  $\Sigma = \mathbb{E}(\mathbf{x}_i \mathbf{x}_i^T) = (\sigma(t_j, t_k))_{1 \leq j \leq n, 1 \leq k \leq n}$  is the covariance matrix of  $X$  at the observations points. Let  $\bar{\mathbf{x}}$  and  $\mathbf{S}$  denote the sample mean and the sample covariance (non corrected by the mean) of the data  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , i.e.

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i, \quad \mathbf{S} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T.$$

Our aim is to build a model selection estimator of the covariance of the process observed with  $N$  replications but without additional assumptions such as stationarity nor Gaussianity. The asymptotics will be taken with respect to  $N$ , the number of copies of the process.

## 2.1 Notations and preliminary definitions

First, define specific matricial notations. We refer to [Lüt96] or [KvR05] for definitions and properties of matrix operations and special matrices. As usual, vectors in  $\mathbb{R}^k$  are regarded as column vectors for all  $k \in \mathbb{N}$ . To be able to write general methods for all our models, we will treat matricial data as a natural extension of the vectorial data, with of course, different correlation structure. For this, we introduce a natural linear transformation, which converts any matrix into a column vector. The vectorization of a  $k \times n$  matrix  $\mathbf{A} = (a_{ij})_{1 \leq i \leq k, 1 \leq j \leq n}$  is the  $kn \times 1$  column vector denoted by  $vec(\mathbf{A})$ , obtain by stacking the columns of the matrix  $\mathbf{A}$  on top of one another. That is  $vec(\mathbf{A}) = [a_{11}, \dots, a_{k1}, a_{12}, \dots, a_{k2}, \dots, a_{1n}, \dots, a_{kn}]^T$ .

For a symmetric  $k \times k$  matrix  $\mathbf{A}$ , the vector  $vec(\mathbf{A})$  contains more information than necessary, since the matrix is completely determined by the lower triangular portion, that is, the  $k(k+1)/2$  entries on and below the main diagonal. Hence, we introduce the symmetrized vectorization, which corresponds to a half-vectorization, denoted by  $vech(\mathbf{A})$ . More precisely, for any matrix  $\mathbf{A} = (a_{ij})_{1 \leq i \leq k, 1 \leq j \leq k}$ , define  $vech(\mathbf{A})$  as the  $k(k+1)/2 \times 1$  column vector obtained by vectorizing only the lower triangular part of  $\mathbf{A}$ . That is  $vech(\mathbf{A}) = [a_{11}, \dots, a_{k1}, a_{22}, \dots, a_{n2}, \dots, a_{(k-1)(k-1)}, a_{(k-1)k}, a_{kk}]^T$ . There exist unique linear transformation which transforms the half-vectorization of a matrix to its vectorization and vice-versa called, respectively, the duplication matrix and the elimination matrix. For any  $k \in \mathbb{N}$ , the  $k^2 \times k(k+1)/2$  duplication matrix is denoted by  $\mathbf{D}_k$ ,  $\mathbf{1}_k = (1, \dots, 1)^T \in \mathbb{R}^k$  and  $\mathbf{I}_k$  is the identity matrix in  $\mathbb{R}^{k \times k}$ .

For any matrix  $\mathbf{A}$ ,  $\mathbf{A}^T$  is the transpose of  $\mathbf{A}$ ,  $tr(\mathbf{A})$  is the trace of  $\mathbf{A}$ ,  $\|\mathbf{A}\|$  is the Frobenius matrix norm defined as  $\|\mathbf{A}\|^2 = tr(\mathbf{A}\mathbf{A}^T)$ ,  $\lambda_{\max}(\mathbf{A})$  is the maximum eigenvalue of  $\mathbf{A}$ ,  $\rho(\mathbf{A})$  is the spectral norm of  $\mathbf{A}$ , that is  $\rho(\mathbf{A}) = \lambda_{\max}(\mathbf{A})$  for  $\mathbf{A}$  a d.n.n matrix. If  $\mathbf{A} = (a_{ij})_{1 \leq i \leq k, 1 \leq j \leq n}$  is a  $k \times n$  matrix and  $\mathbf{B} = (b_{ij})_{1 \leq i \leq p, 1 \leq j \leq q}$  is a  $p \times q$  matrix, then the Kronecker product of the two matrices, denoted by  $\mathbf{A} \otimes \mathbf{B}$ , is the  $kp \times nq$  block matrix

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & \cdot & \cdot & \cdot & a_{1n}\mathbf{B} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ a_{k1}\mathbf{B} & \cdot & \cdot & \cdot & a_{kn}\mathbf{B} \end{bmatrix}.$$

For any random matrix  $\mathbf{Z} = (Z_{ij})_{1 \leq i \leq k, 1 \leq j \leq n}$ , its expectation is denoted by  $\mathbb{E}(\mathbf{Z}) = (\mathbb{E}(Z_{ij}))_{1 \leq i \leq k, 1 \leq j \leq n}$ . For any random vector  $\mathbf{z} = (Z_i)_{1 \leq i \leq k}$ , let  $V(\mathbf{z}) = (cov(Z_i, Z_j))_{1 \leq i, j \leq k}$  be its covariance matrix. With this notation,  $V(\mathbf{x}_1) = V(\mathbf{x}_i) = (\sigma(t_j, t_k))_{1 \leq j \leq n, 1 \leq k \leq n}$  is the covariance matrix of  $X$ .

Let  $(m) \in \mathcal{M}$ , and recall that to the finite set  $\mathcal{G}_m = \{g_\lambda\}_{\lambda \in (m)}$  of functions  $g_\lambda : T \rightarrow \mathbb{R}$  we associate the  $n \times m$  matrix  $\mathbf{G}$  with entries  $g_{j\lambda} = g_\lambda(t_j)$ ,  $j = 1, \dots, n$ ,  $\lambda \in (m)$ . Furthermore, for each  $t \in T$ , we write  $\mathbf{G}_t = (g_\lambda(t), \lambda \in (m))^T$ . For  $k \in \mathbb{N}$ ,  $\mathcal{S}_k$  denotes the linear subspace of  $\mathbb{R}^{k \times k}$  composed of symmetric matrices. For  $\mathbf{G} \in \mathbb{R}^{n \times m}$ ,  $\mathcal{S}(\mathbf{G})$  is the linear subspace of  $\mathbb{R}^{n \times n}$  defined by

$$\mathcal{S}(\mathbf{G}) = \{\mathbf{G}\Psi\mathbf{G}^T : \Psi \in \mathcal{S}_m\}.$$

Let  $\mathcal{S}_N(\mathbf{G})$  be the linear subspace of  $\mathbb{R}^{n \times n}$  defined by

$$\mathcal{S}_N(\mathbf{G}) = \{\mathbf{1}_N \otimes \mathbf{G}\Psi\mathbf{G}^T : \Psi \in \mathcal{S}_m\} = \{\mathbf{1}_N \otimes \Gamma : \Gamma \in \mathcal{S}(\mathbf{G})\}$$

and let  $\mathcal{V}_N(\mathbf{G})$  be the linear subspace of  $\mathbb{R}^{n^2N}$  defined by

$$\mathcal{V}_N(\mathbf{G}) = \{\mathbf{1}_N \otimes \text{vec}(\mathbf{G}\Psi\mathbf{G}^T) : \Psi \in \mathcal{S}_m\} = \{\mathbf{1}_N \otimes \text{vec}(\Gamma) : \Gamma \in \mathcal{S}(\mathbf{G})\}.$$

All these spaces are regarded as Euclidean spaces with the scalar product associated to the Frobenius matrix norm.

## 2.2 Model

The approach that we will develop to estimate the covariance function  $\sigma$  is based on the following two main ingredients: first, we consider a functional expansion  $\tilde{X}$  to approximate the underlying process  $X$  and take the covariance of  $\tilde{X}$  as an approximation of the true covariance  $\Sigma$ .

For this, let  $(m) \in \mathcal{M}$  and consider an approximation to the process  $X$  of the following form:

$$\tilde{X}(t) = \sum_{\lambda \in (m)} a_\lambda g_\lambda(t), \quad (2.1)$$

where  $a_\lambda$  are suitable random coefficients. For instance if  $X$  takes its values in  $L^2(T)$  (the space of square integrable real-valued functions on  $T$ ) and if  $(g_\lambda)_{\lambda \in \mathcal{M}}$  are orthonormal functions in  $L^2(T)$ , then one can take

$$a_\lambda = \int_T X(t)g_\lambda(t)dt.$$

Several basis can thus be considered, such as a polynomial basis on  $\mathbb{R}^d$ , Fourier expansion on a rectangle  $T \subset \mathbb{R}^d$  (i.e.  $g_\lambda(t) = e^{i2\pi(\omega_\lambda \cdot t)}$ , using a regular grid of discrete set of frequencies  $\{\omega_\lambda \in \mathbb{R}^d, \lambda \in (m)\}$  that do not depend on  $t_1, \dots, t_n$ ). One can also use, as in [ETA03], tensorial product of B-splines on a rectangle  $T \subset \mathbb{R}^d$ , with a regular grid of nodes in  $\mathbb{R}^d$  not depending on  $t_1, \dots, t_n$  or a standard wavelet basis on  $\mathbb{R}^d$ , depending on a regular grid of locations in  $\mathbb{R}^d$  and discrete scales in  $\mathbb{R}_+$ . Another class of natural expansion is provided by Karhunen-Loeve expansion of the process  $X$  (see [Adl90] for more references).

Therefore, it is natural to consider the covariance function  $\rho$  of  $\tilde{X}$  as an approximation of  $\sigma$ . Since the covariance  $\rho$  can be written as

$$\rho(s, t) = \mathbf{G}_s^T \overline{\Psi} \mathbf{G}_t, \quad (2.2)$$

where, after reindexing the functions if necessary,  $\mathbf{G}_t = (g_\lambda(t), \lambda \in (m))^T$  and

$$\overline{\Psi} = (\mathbb{E}(a_\lambda a_\mu)), \text{ with } (\lambda, \mu) \in (m) \times (m).$$

Hence we are led to look for an estimate  $\hat{\sigma}$  of  $\sigma$  in the class of functions of the form (2.2), with  $\Psi \in \mathbb{R}^{m \times m}$  some symmetric matrix. Note that the choice of the function expansion in (2.1), in particular the choice the subset of indices  $(m)$ , will be crucial in the approximation properties of the covariance function  $\rho$ . This estimation procedure has several advantages: it will be shown that an appropriate choice of loss function leads to the construction of symmetric d.n.n matrix  $\hat{\Psi}$  (see Proposition 3.1) and thus the resulting estimate

$$\hat{\sigma}(s, t) = \mathbf{G}_s^T \hat{\Psi} \mathbf{G}_t,$$

is a covariance function, so the resulting estimator can be plugged in other procedures which requires working with a covariance function. We also point out that the large amount of existing approaches for function approximation of the type (2.1) (such as those based on Fourier, wavelets, kernel, splines or radial functions) provides great flexibility to the model (2.2).

Secondly, we use the Frobenius matrix norm to quantify the risk of the covariance matrix estimators. Recall that  $\Sigma = (\sigma(t_j, t_k))_{1 \leq j, k \leq n}$  is the true covariance while  $\Gamma = (\rho(t_j, t_k))_{(j, k)}$  will denote be the covariance matrix of the approximated process  $\tilde{X}$  at the observation points. Hence

$$\Gamma = \mathbf{G}\bar{\Psi}\mathbf{G}^T. \quad (2.3)$$

Comparing the covariance function  $\rho$  with the true one  $\sigma$  over the design points  $t_j$ , implies quantifying the deviation of  $\Gamma$  from  $\Sigma$ . For this consider the following loss function

$$L(\Psi) = \mathbb{E} \|\mathbf{x}\mathbf{x}^T - \mathbf{G}\Psi\mathbf{G}^T\|^2,$$

where  $\mathbf{x} = (X(t_1), \dots, X(t_n))^T$  and  $\|\cdot\|$  is the Frobenius matrix norm. Note that

$$L(\Psi) = \|\Sigma - \mathbf{G}\Psi\mathbf{G}^T\|^2 + C,$$

where the constant  $C$  does not depend on  $\Psi$ . The Frobenius matrix norm provides a meaningful metric for comparing covariance matrices, widely used in multivariate analysis, in particular in the theory on principal components analysis. See also [BR97], [SS05] and references therein for other applications of this loss function.

To the loss  $L$  corresponds the following empirical contrast function  $L_N$ , which will be the fitting criterion we will try to minimize

$$L_N(\Psi) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i\mathbf{x}_i^T - \mathbf{G}\Psi\mathbf{G}^T\|^2.$$

We point out that this loss is exactly the sum of the squares of the residuals corresponding to the matrix linear regression model

$$\mathbf{x}_i\mathbf{x}_i^T = \mathbf{G}\Psi\mathbf{G}^T + \mathbf{U}_i, \quad i = 1, \dots, N, \quad (2.4)$$

with i.i.d. matrix errors  $\mathbf{U}_i$  such that  $\mathbb{E}(\mathbf{U}_i) = \mathbf{0}$ . This remark provides a natural framework to study the covariance estimation problem as a matricial regression model. Note also that the set of matrices  $\mathbf{G}\Psi\mathbf{G}^T$  is a linear subspace of  $\mathbb{R}^{n \times n}$  when  $\Psi$  ranges over the space of symmetric matrices  $\mathcal{S}_m$ .

To summarize our approach, we finally propose following two-step estimation procedure: in a first step, for a given design matrix  $\mathbf{G}$ , define

$$\hat{\Psi} = \arg \min_{\Psi \in \mathcal{S}_m} L_N(\Psi),$$

and take  $\hat{\Sigma} = \mathbf{G}\hat{\Psi}\mathbf{G}^T$  as an estimator of  $\Sigma$ . Note that  $\hat{\Psi}$  will be shown to be a d.n.n matrix (see Proposition 3.1) and thus  $\hat{\Sigma}$  is also a d.n.n matrix. Since the minimization of  $L_N(\Psi)$  with respect to  $\Psi$  is done over the linear space of symmetric matrices  $\mathcal{S}_m$ , it can be transformed to a classical least squares linear problem, and the computation of  $\hat{\Psi}$  is therefore quite simple. For a given design matrix  $\mathbf{G}$ , we will construct an estimator

for  $\Gamma = \mathbf{G}\bar{\Psi}\mathbf{G}^T$  which will be close to  $\Sigma = V(\mathbf{x}_1)$  as soon as  $\tilde{X}$  is a sharp estimation of  $X$ . So, the role of  $\mathbf{G}$  and thus the choice of the subset of indices ( $m$ ) is crucial since it determines the behavior of the estimator.

Hence, in second step, we aim at selecting the best design matrix  $\mathbf{G} = \mathbf{G}_m$  among a collection of candidates  $\{\mathbf{G}_m, (m) \in \mathcal{M}\}$ . For this, methods and results from the theory of model selection in linear regression can be applied to the present context. In particular the results in [Bar00], [Com01] or [LL08] will be useful in dealing with model selection for the framework (2.4). Note that only assumptions about moments, not specific distributions of the data, are involved in the estimation procedure.

**Remark 2.1.** We consider here a least-squares estimates of the covariance. Note that suitable regularization terms or constraints could also be incorporated into the minimization of  $L_N(\Psi)$  to impose desired properties for the resulting estimator, such as smoothness or sparsity conditions as in [LRZ08].

### 3 Oracle inequality for Covariance Estimation

The first part of this section describes the properties of the least squares estimator  $\hat{\Sigma} = \mathbf{G}\hat{\Psi}\mathbf{G}^T$  while the second part builds a selection procedure to pick automatically the best estimate among a collection of candidates.

#### 3.1 Least Squares Covariance Estimation

Given some  $n \times m$  fixed design matrix  $\mathbf{G}$  associated to a finite family of  $m$  basis functions, the least squares covariance estimator of  $\Sigma$  is defined by

$$\hat{\Sigma} = \mathbf{G}\hat{\Psi}\mathbf{G}^T = \arg \min \left\{ \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i \mathbf{x}_i^T - \Gamma\|^2 : \Gamma = \mathbf{G}\Psi\mathbf{G}^T, \Psi \in \mathcal{S}_m \right\}. \quad (3.1)$$

The corresponding estimator of the covariance function  $\sigma$  is

$$\hat{\sigma}(s, t) = \mathbf{G}_s^T \hat{\Psi} \mathbf{G}_t. \quad (3.2)$$

**Proposition 3.1.** *Let  $\mathbf{Y}_1, \dots, \mathbf{Y}_N \in \mathbb{R}^{n \times n}$  and  $\mathbf{G} \in \mathbb{R}^{n \times m}$  be arbitrary matrices Then, the infimum*

$$\inf \left\{ \frac{1}{N} \sum_{i=1}^N \|\mathbf{Y}_i - \mathbf{G}\Psi\mathbf{G}^T\|^2 : \Psi \in \mathcal{S}_m \right\}$$

*is achieved at*

$$\hat{\Psi} = (\mathbf{G}^T \mathbf{G})^- \mathbf{G}^T \left( \frac{\bar{\mathbf{Y}} + \bar{\mathbf{Y}}^T}{2} \right) \mathbf{G} (\mathbf{G}^T \mathbf{G})^-, \quad (3.3)$$

*where  $(\mathbf{G}^T \mathbf{G})^-$  is any generalized inverse of  $\mathbf{G}^T \mathbf{G}$  (see [EHN96] for a general definition), and*

$$\bar{\mathbf{Y}} = \frac{1}{N} \sum_{i=1}^N \mathbf{Y}_i.$$

*Furthermore,  $\mathbf{G}\hat{\Psi}\mathbf{G}^T$  is the same for all the generalized inverses  $(\mathbf{G}^T \mathbf{G})^-$  of  $\mathbf{G}^T \mathbf{G}$ . In particular, if  $\mathbf{Y}_1, \dots, \mathbf{Y}_N \in \mathcal{S}_n$  (i.e., if they are symmetric matrices) then any minimizer has the form*

$$\hat{\Psi} = (\mathbf{G}^T \mathbf{G})^- \mathbf{G}^T \bar{\mathbf{Y}} \mathbf{G} (\mathbf{G}^T \mathbf{G})^-.$$



If  $\mathbf{Y}_1, \dots, \mathbf{Y}_N$  are d.n.n. then these matrices  $\widehat{\Psi}$  are d.n.n.

If we assume that  $(\mathbf{G}^T \mathbf{G})^{-1}$  exists, then Proposition 3.1 shows that we retrieve the expression (1.3) for  $\widehat{\Psi}$  that has been derived from least square estimation in model (1.1).

**Theorem 3.2.** Let  $\mathbf{S} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T$ . Then, the least squares covariance estimate defined by (3.1) is given by the d.n.n. matrix

$$\widehat{\Sigma} = \mathbf{G} \widehat{\Psi} \mathbf{G}^T = \Pi \mathbf{S} \Pi,$$

where

$$\begin{aligned} \widehat{\Psi} &= (\mathbf{G}^T \mathbf{G})^{-1} \mathbf{G}^T \mathbf{S} \mathbf{G} (\mathbf{G}^T \mathbf{G})^{-1}, \\ \Pi &= \mathbf{G} (\mathbf{G}^T \mathbf{G})^{-1} \mathbf{G}^T. \end{aligned} \quad (3.4)$$

Moreover  $\widehat{\Sigma}$  has the following interpretations in terms of orthogonal projections:

- i)  $\widehat{\Sigma}$  is the projection of  $\mathbf{S} \in \mathbb{R}^{n \times n}$  on  $\mathcal{S}(\mathbf{G})$ .
- ii)  $\mathbf{1}_N \otimes \widehat{\Sigma}$  is the projection of  $\mathbf{Y} = (\mathbf{x}_1 \mathbf{x}_1^T, \dots, \mathbf{x}_N \mathbf{x}_N^T)^T \in \mathbb{R}^{nN \times n}$  on  $\mathcal{S}_N(\mathbf{G})$ .
- iii)  $\mathbf{1}_N \otimes \text{vec}(\widehat{\Sigma})$  is the projection of  $\mathbf{y} = (\text{vec}^T(\mathbf{x}_1 \mathbf{x}_1^T), \dots, \text{vec}^T(\mathbf{x}_N \mathbf{x}_N^T))^T \in \mathbb{R}^{n^2 N}$  on  $\mathcal{V}_N(\mathbf{G})$ .

The proof of this theorem is a direct application of Proposition 3.1. Hence for a given design matrix  $\mathbf{G}$ , the least squares estimator  $\widehat{\Sigma} = \widehat{\Sigma}(\mathbf{G})$  is well defined and has the structure of a covariance matrix. It remains to study how to pick automatically the estimate when dealing with a collection of design matrices coming from several approximation choices for the random process  $X$ .

## 3.2 Main Result

Consider a collection of indices  $(m) \in \mathcal{M}$  with size  $m$ . Let also  $\{\mathbf{G}_m : (m) \in \mathcal{M}\}$  be a finite family of design matrices  $\mathbf{G}_m \in \mathbb{R}^{n \times m}$ , and let  $\widehat{\Sigma}_m = \widehat{\Sigma}(\mathbf{G}_m)$ ,  $(m) \in \mathcal{M}$ , be the corresponding least squares covariance estimators. The problem of interest is to select the best of these estimators in the sense of the minimal quadratic risk  $\mathbb{E} \left\| \Sigma - \widehat{\Sigma}_m \right\|^2$ .

The main theorem of this section provides a non-asymptotic bound for the risk of a penalized strategy for this problem. For all  $(m) \in \mathcal{M}$ , write

$$\begin{aligned} \Pi_m &= \mathbf{G}_m (\mathbf{G}_m^T \mathbf{G}_m)^{-1} \mathbf{G}_m^T, \\ D_m &= \text{Tr}(\Pi_m), \end{aligned} \quad (3.5)$$

We assume that  $D_m \geq 1$  for all  $(m) \in \mathcal{M}$ . The estimation error for a given model  $(m) \in \mathcal{M}$  is given by

$$\mathbb{E} \left( \left\| \Sigma - \widehat{\Sigma}_m \right\|^2 \right) = \left\| \Sigma - \Pi_m \Sigma \Pi_m \right\|^2 + \frac{\delta_m^2 D_m}{N}, \quad (3.6)$$

where

$$\begin{aligned} \delta_m^2 &= \frac{\text{Tr}((\Pi_m \otimes \Pi_m) \Phi)}{D_m}, \\ \Phi &= V(\text{vec}(\mathbf{x}_1 \mathbf{x}_1^T)). \end{aligned}$$

Given  $\theta > 0$ , define the penalized covariance estimator  $\tilde{\Sigma} = \hat{\Sigma}_{\hat{m}}$  by

$$\hat{m} = \arg \min_{(m) \in \mathcal{M}} \left\{ \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}_i \mathbf{x}_i^T - \hat{\Sigma}_m \right\|^2 + \text{pen}(m) \right\},$$

where

$$\text{pen}(m) = (1 + \theta) \frac{\delta_m^2 D_m}{N}. \quad (3.7)$$

**Theorem 3.3.** *Let  $q > 0$  be given such that there exists  $p > 2(1 + q)$  satisfying  $\mathbb{E} \left\| \mathbf{x}_1 \mathbf{x}_1^T \right\|^p < \infty$ . Then, for some constants  $K(\theta) > 1$  and  $C'(\theta, p, q) > 0$  we have that*

$$\left( \mathbb{E} \left\| \Sigma - \tilde{\Sigma} \right\|^{2q} \right)^{1/q} \leq 2^{(q^{-1}-1)_+} \left[ K(\theta) \inf_{(m) \in \mathcal{M}} \left( \left\| \Sigma - \Pi_m \Sigma \Pi_m \right\|^2 + \frac{\delta_m^2 D_m}{N} \right) + \frac{\Delta_p}{N} \delta_{\text{sup}}^2 \right],$$

where

$$\Delta_p^q = C'(\theta, p, q) \mathbb{E} \left\| \mathbf{x}_1 \mathbf{x}_1^T \right\|^p \left( \sum_{(m) \in \mathcal{M}} \delta_m^{-p} D_m^{-(p/2-1-q)} \right)$$

and

$$\delta_{\text{sup}}^2 = \max \{ \delta_m^2 : (m) \in \mathcal{M} \}.$$

In particular, for  $q = 1$  we have

$$\mathbb{E} \left( \left\| \Sigma - \tilde{\Sigma} \right\|^2 \right) \leq K(\theta) \inf_{(m) \in \mathcal{M}} \mathbb{E} \left( \left\| \Sigma - \hat{\Sigma}_m \right\|^2 \right) + \frac{\Delta_p}{N} \delta_{\text{sup}}^2. \quad (3.8)$$

For the proof of this result, we first restate this theorem in a vectorized form which turns to be a  $d$ -variate extensions of results in [Bar00] (which are covered when  $d = 1$ ) and are stated in Section 4.1. Their proof rely on model selection techniques and a concentration tool stated in Section 4.2.

**Remark 3.4.** The penalty depends on the quantity  $\delta_m$ . Note that

$$\begin{aligned} D_m \delta_m^2 &= \gamma_m^2 = \gamma^2(m, n) = \text{Tr} \left( (\Pi_m \otimes \Pi_m) \Phi \right) \\ &= \mathbb{E} \left\| \hat{\Sigma}_m - \Pi_m \Sigma \Pi_m \right\|^2 N = \text{Tr} \left( V \left( \text{vec} \left( \hat{\Sigma}_m \right) \right) \right) N. \end{aligned} \quad (3.9)$$

So, we get that  $\delta_m^2 \leq \lambda_{\max}(\Phi)$  for all  $(m)$ . Hence Theorem 3.3 remains true if  $\delta_m^2$  is replaced by  $\lambda^2 = \lambda_{\max}(\Phi)$  in all the statements.

**Remark 3.5.** The penalty relies thus on  $\Phi = V(\text{vec}(\mathbf{x}_1 \mathbf{x}_1^T))$ . This quantity reflects the correlation structure of the data. We point out that for practical purpose, this quantity can be estimated using the empirical version of  $\Phi$  since the  $\mathbf{x}_i$ ,  $i = 1, \dots, N$  are i.i.d observed random variables. In the original paper by Baraud [Bar02], an estimator of the variance is proposed to overcome this issue. However, the consistency proof relies on a concentration inequality which turns to be a  $\chi^2$  like inequality. Extending this inequality to our case would mean to be able to construct concentration bounds for matrices  $\mathbf{x} \mathbf{x}^T$ , implying Wishart distributions. If some results exist in this framework [RMSE08], adapting this kind of construction to our case falls beyond the scope of this paper.

We have obtained in Theorem 3.3 an oracle inequality since, using (3.6) and (3.8), one immediately sees that  $\tilde{\Sigma}$  has the same quadratic risk as the “oracle” estimator except for an additive term of order  $O\left(\frac{1}{N}\right)$  and a constant factor. Hence, the selection procedure is optimal in the sense that it behaves as if the true model were at hand. To describe the result in terms of rate of convergence, we have to pay a special attention to the bias terms  $\|\Sigma - \Pi_m \Sigma \Pi_m\|^2$ . In a very general framework, it is difficult to evaluate such approximation terms. If the process has bounded second moments, i.e for all  $i = 1, \dots, n$ , we have  $\mathbb{E}(X^2(t_i)) \leq C$ , then we can write

$$\begin{aligned} \|\Sigma - \Pi_m \Sigma \Pi_m\|^2 &\leq C_2 \sum_{i=1}^n \sum_{i'=1}^n \left[ \mathbb{E} \left( X(t_i) - \tilde{X}(t_i) \right)^2 + \mathbb{E} \left( X(t_{i'}) - \tilde{X}(t_{i'}) \right)^2 \right] \\ &\leq 2C_2 n^2 \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left( X(t_i) - \tilde{X}(t_i) \right)^2. \end{aligned}$$

Since  $n$  is fixed and the asymptotics are given with respect to  $N$ , the number of replications of the process, the rate of convergence relies on the quadratic error of the expansion of the process.

For example take  $d = 1$ ,  $T = [a, b]$ ,  $\mathcal{M} = \mathcal{M}_N = \{(m) = \{1, \dots, m\}, m = 1, \dots, N\}$ , and for a process  $X(t)$  with  $t \in [a, b]$ , consider its Karhunen-Loève expansion (see for instance [Adl90]), i.e. write

$$X(t) = \sum_{\lambda=1}^{\infty} Z_{\lambda} g_{\lambda}(t),$$

where  $Z_{\lambda}$  are centered random variables with  $\mathbb{E}(Z_{\lambda}^2) = \gamma_{\lambda}^2$ , where  $\gamma_{\lambda}^2$  is the eigenvalue corresponding to the eigenfunction  $g_{\lambda}$  of the operator  $(Kf)(t) = \int_a^b \sigma(s, t) f(s) ds$ . If  $X(t)$  is a Gaussian process then the random variables  $Z_{\lambda}$  are Gaussian and stochastically independent. Hence, a natural approximation of  $X(t)$  is given by

$$\tilde{X}(t) = \sum_{\lambda=1}^m Z_{\lambda} g_{\lambda}(t).$$

So we have that

$$\mathbb{E} \left( X(t) - \tilde{X}(t) \right)^2 = \mathbb{E} \left( \sum_{\lambda=m+1}^{\infty} Z_{\lambda} g_{\lambda}(t) \right)^2 = \sum_{\lambda=m+1}^{\infty} \gamma_{\lambda}^2 g_{\lambda}^2(t).$$

therefore, if  $\|g_{\lambda}\|_{L_2([a,b])}^2 = 1$  then  $\mathbb{E} \left\| X(t) - \tilde{X}(t) \right\|_{L_2([a,b])}^2 = \sum_{l=m+1}^{\infty} \gamma_{\lambda}^2$ . Assume that the  $\gamma_{\lambda}$ 's have a polynomial decay of rate  $\alpha > 0$ , namely  $\gamma_{\lambda} \sim \lambda^{-\alpha}$ , then we get an approximation error of order  $O\left((m+1)^{-2\alpha}\right)$ . Hence, we get that (under appropriate conditions on the design points  $t_1, \dots, t_n$ )

$$\|\Sigma - \Pi_m \Sigma \Pi_m\|^2 = O\left((m+1)^{-2\alpha}\right).$$

Finally, since in this example  $\mathbb{E} \left\| \Sigma - \tilde{\Sigma} \right\|^2 \leq K(\theta) \inf_{m \in \mathcal{M}_N} \left( \|\Sigma - \Pi_m \Sigma \Pi_m\|^2 + \frac{\delta_m^2 m}{N} \right) + O\left(\frac{1}{N}\right)$  then the quadratic risk is of order  $N^{-\frac{2\alpha}{2\alpha+1}}$  as soon as  $m \sim N^{1/(2\alpha+1)}$  belongs to the

collection of models  $\mathcal{M}_N$ . In another framework, if we consider a spline expansion, the rate of convergence for the approximation given in [ETA03] are of the same order.

Hence we have obtained a model selection procedure which enables to recover the best covariance model among a given collection. This method works without strong assumptions on the process, in particular stationarity is not assumed, but at the expense of necessary i.i.d observations of the process at the same points. However the range of applications is broad, especially in geophysics or epidemiology.

## 4 Model Selection for Multidimensional Regression

### 4.1 Oracle Inequality for multidimensional regression model

Recall that we consider the following model

$$\mathbf{x}_i \mathbf{x}_i^T = \mathbf{G} \Psi \mathbf{G}^T + \mathbf{U}_i, \quad i = 1, \dots, N,$$

with i.i.d. matrix errors  $\mathbf{U}_i$ ,  $\mathbb{E}(\mathbf{U}_i) = \mathbf{0}$ . This model can be equivalently rewritten in vectorized form in the following way

$$\mathbf{y} = \mathbf{A} \beta + \mathbf{u},$$

where  $\mathbf{y}$  is a data vector,  $\mathbb{E}(\mathbf{u}) = \mathbf{0}$ ,  $\mathbf{A}$  is a known fixed matrix, and  $\beta = \text{vech}(\Psi)$  is an unknown vector parameter. It is worth of noting that this regression model has several peculiarities in comparison with standard ones.

*i)* The error  $\mathbf{u}$  has a specific correlation structure, namely  $\mathbf{I}_N \otimes \Phi$ , where  $\Phi = V(\text{vec}(\mathbf{x}_i \mathbf{x}_i^T))$ .

*ii)* In contrast with standard multivariate models, each coordinate of  $\mathbf{y}$  depends on all the coordinates of  $\beta$ .

*iii)* For any estimator  $\hat{\Sigma} = \mathbf{G} \hat{\Psi} \mathbf{G}^T$  that be a linear function of the sample covariance  $\mathbf{S}$  of the data  $\mathbf{x}_1, \dots, \mathbf{x}_N$  (and so, in particular, for the estimator minimizing  $L_N$ ) it is possible to construct an unbiased estimator of its quadratic risk  $\mathbb{E} \left\| \Sigma - \hat{\Sigma} \right\|^2$ .

Assume we observe  $\mathbf{y}_i$ ,  $i = 1, \dots, N$  random vectors of  $\mathbb{R}^d$  such that

$$\mathbf{y}_i = \mathbf{f}^i + \varepsilon_i, \quad i = 1, \dots, N, \quad (4.1)$$

where  $\mathbf{f}^i \in \mathbb{R}^d$  are nonrandom and  $\varepsilon_1, \dots, \varepsilon_N$  are i.i.d. random vectors in  $\mathbb{R}^d$  with  $E(\varepsilon_1) = \mathbf{0}$  and  $V(\varepsilon_1) = \Phi$ . For sake of simplicity, we identify the function  $g : \mathcal{X} \rightarrow \mathbb{R}^d$  with vectors  $(g(x_1) \dots g(x_N))^T \in \mathbb{R}^{Nd}$  and we denote by  $\langle a, b \rangle_N = \frac{1}{N} \sum_{i=1}^N a_i^T b_i$ , with  $a = (a_1 \dots a_N)^T$  and  $a_i \in \mathbb{R}^d$ , the inner product of  $\mathbb{R}^{Nd}$  associated to the norm  $\|\cdot\|_N$ .

Given  $N, d \in \mathbb{N}$ , let  $(\mathcal{L}_m)_{(m) \in \mathcal{M}}$  be a finite family of linear subspaces of  $\mathbb{R}^{Nd}$ . For each  $(m) \in \mathcal{M}$ , assume  $\mathcal{L}_m$  has dimension  $D_m \geq 1$ . For each  $(m) \in \mathcal{M}$ , let  $\hat{\mathbf{f}}_m$  be the least squares estimator of  $\mathbf{f} = \left( (\mathbf{f}^1)^T, \dots, (\mathbf{f}^N)^T \right)^T$  based on the data  $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_N)$  under the model  $\mathcal{L}_m$ ; i.e.,

$$\hat{\mathbf{f}}_m = \arg \min_{\mathbf{v} \in \mathcal{L}_m} \{ \|\mathbf{y} - \mathbf{v}\|_N^2 \} = \mathbf{P}_m \mathbf{y},$$

where  $\mathbf{P}_m$  is the projector matrix from  $\mathbb{R}^{Nd}$  on  $\mathcal{L}_m$ . Write

$$\delta_m^2 = \frac{\text{Tr}(\mathbf{P}_m (\mathbf{I}_N \otimes \Phi))}{D_m},$$

$$\delta_{\text{sup}}^2 = \max \{ \delta_m^2 : m \in \mathcal{M} \}.$$

Given  $\theta > 0$ , define the penalized estimator  $\tilde{\mathbf{f}} = \widehat{\mathbf{f}}_{\widehat{m}}$ , where

$$\widehat{m} = \arg \min_{(m) \in \mathcal{M}} \left\{ \left\| \mathbf{y} - \widehat{\mathbf{f}}_m \right\|_N^2 + \text{pen}(m) \right\},$$

with

$$\text{pen}(m) = (1 + \theta) \frac{\delta_m^2 D_m}{N}.$$

**Proposition 4.1.** : *Let  $q > 0$  be given such that there exists  $p > 2(1 + q)$  satisfying  $\mathbb{E} \|\varepsilon_1\|^p < \infty$ . Then, for some constants  $K(\theta) > 1$  and  $c(\theta, p, q) > 0$  we have that*

$$\mathbb{E} \left( \left\| \mathbf{f} - \tilde{\mathbf{f}} \right\|_N^2 - K(\theta) \mathcal{M}^* \right)_+^q \leq \Delta_p^q \frac{\delta_{\text{sup}}^{2q}}{N^q}, \quad (4.2)$$

where

$$\Delta_p^q = C(\theta, p, q) \mathbb{E} \|\varepsilon_1\|^p \left( \sum_{m \in \mathcal{M}} \delta_m^{-p} D_m^{-(p/2-1-q)} \right),$$

$$\mathcal{M}^* = \inf_{(m) \in \mathcal{M}} \left\{ \left\| \mathbf{f} - \mathbf{P}_m \mathbf{f} \right\|_N^2 + \frac{\delta_m^2 D_m}{N} \right\}.$$

This theorem is equivalent to Theorem 3.3 using the vectorized version of the model (4.1) and turns to be an extension of Theorem 3.1 in [Bar00] to the multivariate case. In a similar way, the following result constitutes also a natural extension of Corollary 3.1 in [Bar00]. It is also closely related to the recent work in [Gen08].

**Corollary 4.2.** . *Under the assumptions of Proposition 4.1 it holds that*

$$\left( \mathbb{E} \left\| \mathbf{f} - \tilde{\mathbf{f}} \right\|_N^{2q} \right)^{1/q} \leq 2^{(q-1)_+} \left[ K(\theta) \inf_{m \in \mathcal{M}} \left( \left\| \mathbf{f} - \mathbf{P}_m \mathbf{f} \right\|^2 + \frac{\delta_m^2 D_m}{N} \right) + \frac{\Delta_p}{N} \delta_{\text{sup}}^2 \right],$$

where  $\Delta_p$  was defined in Proposition (4.1).

Under regularity assumptions for the function  $\mathbf{f}$ , depending on a smoothness parameter  $s$ , the bias term is of order

$$\left\| \mathbf{f} - \mathbf{P}_m \mathbf{f} \right\|^2 = O(D_m^{-2s}).$$

Hence, for  $q = 1$  we obtain the usual rate of convergence  $N^{-\frac{2s}{2s+1}}$  for the quadratic risk as soon as the optimal choice  $D_m = N^{\frac{1}{2s+1}}$  belongs to the collection of models, yielding the optimal rate of convergence for the penalized estimator.

## 4.2 Concentration Bound for multidimensional random process

These results are  $d$ -variate extensions of results in [Bar00] (which are covered when  $d = 1$ ). Their proofs are deferred to the Appendix.

**Proposition 4.3.** (*Extension of Corollary 5.1 in [Bar00]*). *Given  $N, d \in \mathbb{N}$ , let  $\tilde{\mathbf{A}} \in \mathbb{R}^{Nd \times Nd} \setminus \{\mathbf{0}\}$  be a n.n.d. matrix and  $\varepsilon_1, \dots, \varepsilon_N$  i.i.d random vectors in  $\mathbb{R}^d$  with  $\mathbb{E}(\varepsilon_1) = 0$  and  $V(\varepsilon_1) = \Phi$ . Write  $\varepsilon = (\varepsilon_1^T, \dots, \varepsilon_N^T)^T$ ,  $\zeta(\varepsilon) = \sqrt{\varepsilon^T \tilde{\mathbf{A}} \varepsilon}$ , and  $\gamma^2 = \text{Tr}(\tilde{\mathbf{A}}(\mathbf{I}_N \otimes \Phi)) = \delta^2 \text{Tr}(\tilde{\mathbf{A}})$ . For all  $p \geq 2$  such that  $\mathbb{E} \|\varepsilon_1\|^p < \infty$  it holds that, for all  $x > 0$*

$$\mathbb{P} \left( \zeta^2(\varepsilon) \geq \delta^2 \text{Tr}(\tilde{\mathbf{A}}) + 2\delta^2 \sqrt{\text{Tr}(\tilde{\mathbf{A}})} \delta x + \delta^2 \text{Tr}(\tilde{\mathbf{A}}) x \right) \leq C(p) \frac{\mathbb{E} \|\varepsilon_1\|^p \text{Tr}(\tilde{\mathbf{A}})}{\delta^p \rho(\tilde{\mathbf{A}}) x^{p/2}}, \quad (4.3)$$

where the constant  $C(p)$  depends only on  $p$ .

Proposition 4.3 reduces to Corollary 5.1 in [Bar00] when when we only consider  $d = 1$ , in which case  $\delta^2 = (\Phi)_{11} = \sigma^2$  is the variance of the univariate i.i.d. errors  $\varepsilon_i$ .

## 5 Appendix

### 5.1 Proofs of Preliminar results

Proof of Proposition 3.1

*Proof.* a) The minimization problem posed in this theorem is equivalent to minimize

$$h(\Psi) = \|\bar{\mathbf{Y}} - \mathbf{G}\Psi\mathbf{G}^T\|^2.$$

The Frobenius norm  $\|\cdot\|$  is invariant by the *vec* operation. Furthermore,  $\Psi \in \mathcal{S}_m$  can be represented by means of  $\delta = \text{vec}(\Psi) = \mathbf{D}_q\beta$  where  $\beta \in \mathbb{R}^{q(q+1)/2}$ . These facts and the identity

$$\text{vec}(\mathbf{ABC}) = (\mathbf{C}^T \otimes \mathbf{A}) \text{vec}(\mathbf{B}) \quad (5.1)$$

allow one to rewrite

$$h(\Psi) = \|\bar{\mathbf{y}} - (\mathbf{G} \otimes \mathbf{G}) \mathbf{D}_q\beta\|^2,$$

where  $\bar{\mathbf{y}} = \text{vec}(\bar{\mathbf{Y}})$ . Minimization of this quadratic function with respect to  $\beta$  in  $\mathbb{R}^{q(q+1)/2}$  is equivalent to solve the normal equation

$$\mathbf{D}_q^T (\mathbf{G} \otimes \mathbf{G})^T (\mathbf{G} \otimes \mathbf{G}) \mathbf{D}_q\beta = \mathbf{D}_q^T (\mathbf{G} \otimes \mathbf{G})^T \bar{\mathbf{y}}.$$

By using the identities

$$\mathbf{D}_q^T \text{vec}(\mathbf{A}) = \text{vech}(\mathbf{A} + \mathbf{A}^T - \text{diag}(\mathbf{A}))$$

and 5.1, said normal equation can be rewritten

$$\text{vech}(\mathbf{G}^T \mathbf{G} (\Psi + \Psi^T) \mathbf{G}^T \mathbf{G} - \text{diag}(\mathbf{G}^T \mathbf{G} \Psi \mathbf{G}^T \mathbf{G})) = \text{vech}\left(\mathbf{G}^T \left(\bar{\mathbf{Y}} + \bar{\mathbf{Y}}^T\right) \mathbf{G}\right).$$

Finally, it can be verified that  $\hat{\Psi}$  given by (3.3) satisfies this equation as a consequence of the fact that such  $\hat{\Psi}$  it holds that

$$\mathbf{G}^T \mathbf{G} \hat{\Psi} \mathbf{G}^T \mathbf{G} = \text{vech}\left(\mathbf{G}^T \left(\frac{\bar{\mathbf{Y}} + \bar{\mathbf{Y}}^T}{2}\right) \mathbf{G}\right).$$

b) It straightforwardly follows from part a). □

### 5.2 Proofs of Main Results

Proof of Proposition (4.1)

*Proof.* The proof follows the guidelines of the proof in [Bar00]. More generally we will prove that for any  $\eta > 0$  and any sequence of positive numbers  $L_m$ , if the penalty function  $\text{pen} : \mathcal{M} \rightarrow \mathbb{R}_+$  is chosen to satisfy:

$$\text{pen}(m) = (1 + \eta + L_m) \frac{\delta_m^2}{N} D_m \text{ for all } (m) \in \mathcal{M}, \quad (5.2)$$

then for each  $x > 0$  and  $p \geq 2$

$$\mathbb{P} \left( \mathcal{H}(\mathbf{f}) \geq \left(1 + \frac{2}{\eta}\right) \frac{x}{N} \delta_m^2 \right) \leq c(p, \eta) \mathbb{E} \|\varepsilon_1\|^p \sum_{(m) \in \mathcal{M}} \frac{1}{\delta_m^p} \frac{D_m \vee 1}{(L_m D_m + x)^{p/2}}, \quad (5.3)$$

where we have set

$$\mathcal{H}(\mathbf{f}) = \left[ \left\| \mathbf{f} - \tilde{\mathbf{f}} \right\|_N^2 - \left(2 - \frac{4}{\eta}\right) \inf_{(m) \in \mathcal{M}} \{d_N^2(\mathbf{f}, \mathcal{L}_m) + \text{pen}(m)\} \right]_+.$$

To obtain (4.2), take  $\eta = \frac{\theta}{2} = L_m$ . As for each  $(m) \in \mathcal{M}$ ,

$$\begin{aligned} d_N^2(\mathbf{f}, \mathcal{L}_m) + \text{pen}(m) &\leq d_N^2(\mathbf{f}, \mathcal{L}_m) + (1 + \theta) \frac{\delta_m^2}{N} D_m \\ &\leq (1 + \theta) \left( d_N^2(\mathbf{f}, \mathcal{L}_m) + \frac{\delta_m^2}{N} D_m \right) \end{aligned}$$

we get that for all  $q > 0$ ,

$$\mathcal{H}^q(\mathbf{f}) \geq \left[ \left\| \mathbf{f} - \tilde{\mathbf{f}} \right\|_N^2 - \left(2 + \frac{8}{\theta}\right) (1 + \theta) \mathcal{M}^* \right]_+^q = \left[ \left\| \mathbf{f} - \tilde{\mathbf{f}} \right\|_N^2 - K(\theta) \mathcal{M}^* \right]_+^q, \quad (5.4)$$

where  $K(\theta) = \left(2 + \frac{8}{\theta}\right) (1 + \theta)$ .

Since

$$\mathbb{E}(\mathcal{H}^q(\mathbf{f})) = \int_0^\infty q u^{q-1} \mathbb{P}(\mathcal{H}(\mathbf{f}) > u) du,$$

we derive from (5.4) and (5.3) that for all  $p > 2(1 + q)$

$$\begin{aligned} \mathbb{E} \left[ \left( \left\| \mathbf{f} - \tilde{\mathbf{f}} \right\|_N^2 - K(\theta) \mathcal{M}^* \right)_+^q \right] &\leq \mathbb{E}(\mathcal{H}^q(\mathbf{f})) \\ &\leq c(p, \theta) \left(1 + \frac{4}{\theta}\right)^q \frac{\mathbb{E} \|\varepsilon_1\|^p}{N^q} \sum_{(m) \in \mathcal{M}} \frac{\delta_m^{2q}}{\delta_m^p} \int_0^\infty q x^{q-1} \left[ \frac{D_m \vee 1}{\left(\frac{\theta}{2} D_m + x\right)^{p/2}} \wedge 1 \right] dx \\ &\leq c'(p, q, \theta) \frac{\mathbb{E} \|\varepsilon_1\|^p}{N^q} \delta_{\text{sup}}^{2q} \left[ \sum_{(m) \in \mathcal{M}} \delta_m^{-p} D_m^{-(p/2-1-q)} \right] \end{aligned}$$

using that  $\mathbb{P}(\mathcal{H}(\mathbf{f}) > u) \leq 1$ .

Indeed, for  $m \in \mathcal{M}$  such that  $D_m \geq 1$ , using that  $q - 1 - p/2 < 0$ , we get the following

bounds

$$\begin{aligned}
\frac{\delta_m^{2q}}{\delta_m^p} \int_0^\infty qx^{q-1} \left[ \frac{D_m \vee 1}{\left(\frac{\theta}{2}D_m + x\right)^{p/2}} \wedge 1 \right] dx &\leq \delta_{\text{sup}}^{2q} \delta_m^{-p} \int_0^\infty qx^{q-1} \left[ \frac{D_m}{\left(\frac{\theta}{2}D_m + x\right)^{p/2}} \right] dx \\
&= \delta_{\text{sup}}^{2q} \delta_m^{-p} \left( \int_0^{D_m} qx^{q-1} \left[ \frac{D_m}{\left(\frac{\theta}{2}D_m + x\right)^{p/2}} \right] dx + \int_{D_m}^\infty qx^{q-1} \left[ \frac{D_m}{\left(\frac{\theta}{2}D_m + x\right)^{p/2}} \right] dx \right) \\
&\leq \delta_{\text{sup}}^{2q} \delta_m^{-p} \left( \frac{D_m}{\left(\frac{\theta}{2}D_m\right)^{p/2}} \int_0^{D_m} qx^{q-1} dx + D_m \int_{D_m}^\infty qx^{q-1} \left[ \frac{1}{x^{p/2}} \right] dx \right) \\
&= \delta_{\text{sup}}^{2q} \delta_m^{-p} \left( 2^{p/2} \theta^{-p/2} D_m^{1-p/2} \int_0^{D_m} qx^{q-1} dx + D_m \int_{D_m}^\infty qx^{q-1-p/2} dx \right) \\
&= \delta_{\text{sup}}^{2q} \delta_m^{-p} \left( 2^{p/2} \theta^{-p/2} D_m^{1-p/2} [D_m^q] + D_m \left[ \frac{q}{p/2 - q} D_m^{q-p/2} \right] \right) \\
&= \delta_{\text{sup}}^{2q} \delta_m^{-p} \left( 2^{p/2} \theta^{-p/2} D_m^{1-p/2+q} + D_m^{1-p/2+q} \left[ \frac{q}{p/2 - q} \right] \right) \\
&= \delta_{\text{sup}}^{2q} \delta_m^{-p} \left( D_m^{-(p/2-1-q)} \left[ 2^{p/2} \theta^{-p/2} + \frac{q}{p/2 - q} \right] \right). \tag{5.5}
\end{aligned}$$

(5.5) enables to conclude that (4.2) holds assuming (5.3).

We now turn to the proof of (5.3). Recall that, we identify the function  $g : \mathcal{X} \rightarrow \mathbb{R}^d$  with vectors  $(g(x_1) \dots g(x_N))^T \in \mathbb{R}^{Nd}$  and we define the empirical scalar product as  $\langle a, b \rangle_N = \frac{1}{N} \sum_{i=1}^N a_i^T b_i$ , with  $a = (a_1 \dots a_N)^T$  and  $a_i \in \mathbb{R}^d$ , the inner product of  $\mathbb{R}^{Nd}$  associated to the norm  $\|\cdot\|_N$ . For each  $(m) \in \mathcal{M}$  we denote by  $\mathbf{P}_m$  the orthogonal projector onto the linear space  $\left\{ (g(x_1) \dots g(x_N))^T : g \in \mathcal{L}_m \right\} \subset \mathbb{R}^{Nd}$ . This linear space is also denoted by  $\mathcal{L}_m$ . From now on, the subscript  $m$  denotes any minimizer of the function  $m' \rightarrow \|\mathbf{f} - \mathbf{P}_{m'} \mathbf{f}\|^2 + \text{pen}(m')$ ,  $(m') \in \mathcal{M}_N$ . For any  $\mathbf{g} \in \mathbb{R}^{Nd}$  we define the least-squares loss function by

$$\gamma_N(\mathbf{g}) = \|\mathbf{y} - \mathbf{g}\|_N^2$$

Using the definition of  $\gamma_N$  we have that for all  $\mathbf{g} \in \mathbb{R}^{Nd}$ ,

$$\gamma_N(\mathbf{g}) = \|\mathbf{f} + \varepsilon - \mathbf{g}\|_N^2.$$

Then we derive that

$$\|\mathbf{f} - \mathbf{g}\|_N^2 = \gamma_N(\mathbf{f}) + 2 \langle \mathbf{f} - \mathbf{y}, \varepsilon \rangle_N + \|\varepsilon\|_N^2$$

and therefore

$$\|\mathbf{f} - \tilde{\mathbf{f}}\|_N^2 - \|\mathbf{f} - \mathbf{P}_m \mathbf{f}\|_N^2 = \gamma_N(\tilde{\mathbf{f}}) - \gamma_N(\mathbf{P}_m \mathbf{f}) + 2 \langle \tilde{\mathbf{f}} - \mathbf{P}_m \mathbf{f}, \varepsilon \rangle_N. \tag{5.6}$$

By the definition of  $\tilde{\mathbf{f}}$ , we know that

$$\gamma_N(\tilde{\mathbf{f}}) + \text{pen}(\hat{m}) \leq \gamma_N(\mathbf{g}) + \text{pen}(m)$$



for all  $(m) \in \mathcal{M}$  and for all  $\mathbf{g} \in \mathcal{L}_m$ . Then

$$\gamma_N(\tilde{\mathbf{f}}) - \gamma_N(\mathbf{P}_m \mathbf{f}) \leq \text{pen}(m) - \text{pen}(\hat{m}). \quad (5.7)$$

So we get from (5.6) and (5.7) that

$$\left\| \mathbf{f} - \tilde{\mathbf{f}} \right\|_N^2 \leq \left\| \mathbf{f} - \mathbf{P}_m \mathbf{f} \right\|_N^2 + \text{pen}(m) - \text{pen}(\hat{m}) + 2 \langle \mathbf{f} - \mathbf{P}_m \mathbf{f}, \varepsilon \rangle_N + 2 \langle \mathbf{P}_{\hat{m}} \mathbf{f} - \mathbf{f}, \varepsilon \rangle_N + 2 \left\langle \tilde{\mathbf{f}} - \mathbf{P}_{\hat{m}} \mathbf{f}, \varepsilon \right\rangle_N. \quad (5.8)$$

In the following we set for each  $(m') \in \mathcal{M}$ ,

$$\begin{aligned} \mathcal{B}_{m'} &= \{ \mathbf{g} \in \mathcal{L}_{m'} : \|\mathbf{g}\|_N \leq 1 \}, \\ G_{m'} &= \sup_{t \in \mathcal{B}_{m'}} \langle \mathbf{g}, \varepsilon \rangle_N = \|\mathbf{P}_{m'} \varepsilon\|_N, \\ \mathbf{u}_{m'} &= \begin{cases} \frac{\mathbf{P}_{m'} \mathbf{f} - \mathbf{f}}{\|\mathbf{P}_{m'} \mathbf{f} - \mathbf{f}\|_N} & \text{if } \|\mathbf{P}_{m'} \mathbf{f} - \mathbf{f}\|_N \neq 0 \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

Since  $\tilde{\mathbf{f}} = \mathbf{P}_{\hat{m}} \mathbf{f} + \mathbf{P}_{\hat{m}} \varepsilon$ , (5.8) gives

$$\begin{aligned} \left\| \mathbf{f} - \tilde{\mathbf{f}} \right\|_N^2 &\leq \left\| \mathbf{f} - \mathbf{P}_m \mathbf{f} \right\|_N^2 + \text{pen}(m) - \text{pen}(\hat{m}) \\ &\quad + 2 \|\mathbf{f} - \mathbf{P}_m \mathbf{f}\|_N |\langle \mathbf{u}_m, \varepsilon \rangle_N| + 2 \|\mathbf{f} - \mathbf{P}_{\hat{m}} \mathbf{f}\|_N |\langle \mathbf{u}_{\hat{m}}, \varepsilon \rangle_N| + 2G_{\hat{m}}^2. \end{aligned} \quad (5.9)$$

Using repeatedly the following elementary inequality that holds for all positive numbers  $\alpha, x, z$

$$2xz \leq \alpha x^2 + \frac{1}{\alpha} z^2 \quad (5.10)$$

we get for any  $m' \in \mathcal{M}$

$$2 \|\mathbf{f} - \mathbf{P}_{m'} \mathbf{f}\| |\langle \mathbf{u}_{m'}, \varepsilon \rangle_N| \leq \alpha \|\mathbf{f} - \mathbf{P}_{m'} \mathbf{f}\|_N^2 + \frac{1}{\alpha} |\langle \mathbf{u}_{m'}, \varepsilon \rangle_N|^2. \quad (5.11)$$

By Pythagoras Theorem we have

$$\begin{aligned} \left\| \mathbf{f} - \tilde{\mathbf{f}} \right\|_N^2 &= \left\| \mathbf{f} - \mathbf{P}_{\hat{m}} \mathbf{f} \right\|_N^2 + \left\| \mathbf{P}_{\hat{m}} \mathbf{f} - \tilde{\mathbf{f}} \right\|_N^2 \\ &= \left\| \mathbf{f} - \mathbf{P}_{\hat{m}} \mathbf{f} \right\|_N^2 + G_{\hat{m}}^2. \end{aligned} \quad (5.12)$$

We derive from (5.9) and (5.11) that for any  $\alpha > 0$ :

$$\begin{aligned} \left\| \mathbf{f} - \tilde{\mathbf{f}} \right\|_N^2 &\leq \left\| \mathbf{f} - \mathbf{P}_m \mathbf{f} \right\|_N^2 + \alpha \|\mathbf{f} - \mathbf{P}_m \mathbf{f}\|_N^2 + \frac{1}{\alpha} \langle \mathbf{u}_m, \varepsilon \rangle_N^2 \\ &\quad + \alpha \|\mathbf{f} - \mathbf{P}_{\hat{m}} \mathbf{f}\|_N^2 + \frac{1}{\alpha} \langle \mathbf{u}_{\hat{m}}, \varepsilon \rangle_N^2 + 2G_{\hat{m}}^2 + \text{pen}(m) - \text{pen}(\hat{m}). \end{aligned}$$

Now taking into account that by equation (5.12)  $\|\mathbf{f} - \mathbf{P}_{\hat{m}} \mathbf{f}\|_N^2 = \left\| \mathbf{f} - \tilde{\mathbf{f}} \right\|_N^2 - G_{\hat{m}}^2$  the above inequality is equivalent to:

$$(1 - \alpha) \left\| \mathbf{f} - \tilde{\mathbf{f}} \right\|_N^2 \leq (1 + \alpha) \|\mathbf{f} - \mathbf{P}_m \mathbf{f}\|_N^2 + \frac{1}{\alpha} \langle \mathbf{u}_m, \varepsilon \rangle_N^2$$

$$+ \frac{1}{\alpha} \langle \mathbf{u}_{\widehat{m}}, \varepsilon \rangle_N^2 + (2 - \alpha) G_{\widehat{m}}^2 + \text{pen}(m) - \text{pen}(\widehat{m}). \quad (5.13)$$

We choose  $\alpha = \frac{2}{2+\eta} \in ]0, 1[$ , but for sake of simplicity we keep using the notation  $\alpha$ . Let  $\tilde{p}_1$  and  $\tilde{p}_2$  be two functions depending on  $\eta$  mapping  $\mathcal{M}$  into  $\mathbb{R}_+$ . They will be specified later to satisfy

$$\text{pen}(m') \geq (2 - \alpha) \tilde{p}_1(m') + \frac{1}{\alpha} \tilde{p}_2(m') \quad \forall (m') \in \mathcal{M}. \quad (5.14)$$

Since  $\frac{1}{\alpha} \tilde{p}_2(m') \leq \text{pen}(m')$  and  $1 + \alpha \leq 2$ , we get from (5.13) and (5.14) that

$$\begin{aligned} (1 - \alpha) \left\| \mathbf{f} - \tilde{\mathbf{f}} \right\|_N^2 &\leq (1 + \alpha) \left\| \mathbf{f} - \mathbf{P}_m \mathbf{f} \right\|_N^2 + \text{pen}(m) + \frac{1}{\alpha} \tilde{p}_2(m) + (2 - \alpha) (G_{\widehat{m}}^2 - \tilde{p}_1(\widehat{m})) \\ &\quad + \frac{1}{\alpha} (\langle \mathbf{u}_{\widehat{m}}, \varepsilon \rangle_N^2 - \tilde{p}_2(\widehat{m})) + \frac{1}{\alpha} (\langle \mathbf{u}_m, \varepsilon \rangle_N^2 - \tilde{p}_2(m)) \\ &\leq 2 \left( \left\| \mathbf{f} - \mathbf{P}_m \mathbf{f} \right\|_N^2 + \text{pen}(m) \right) + (2 - \alpha) (G_{\widehat{m}}^2 - \tilde{p}_1(\widehat{m})) \\ &\quad + \frac{1}{\alpha} (\langle \mathbf{u}_{\widehat{m}}, \varepsilon \rangle_N^2 - \tilde{p}_2(\widehat{m})) + \frac{1}{\alpha} (\langle \mathbf{u}_m, \varepsilon \rangle_N^2 - \tilde{p}_2(m)). \end{aligned} \quad (5.15)$$

As  $\frac{2}{1-\alpha} = 2 + \frac{4}{\eta}$  we obtain that

$$\begin{aligned} (1 - \alpha) \mathcal{H}(\mathbf{f}) &= \left\{ (1 - \alpha) \left\| \mathbf{f} - \tilde{\mathbf{f}} \right\|_N^2 - (1 - \alpha) \left( 2 + \frac{4}{\eta} \right) \inf_{m' \in \mathcal{M}} \left( \left\| \mathbf{f} - \mathbf{P}_{m'} \mathbf{f} \right\|_N^2 + \text{pen}(m') \right) \right\}_+ \\ &= \left\{ (1 - \alpha) \left\| \mathbf{f} - \tilde{\mathbf{f}} \right\|_N^2 - 2 \left( \left\| \mathbf{f} - \mathbf{P}_m \mathbf{f} \right\|_N^2 + 2 \text{pen}(m) \right) \right\}_+ \\ &\leq \left\{ (2 - \alpha) (G_{\widehat{m}}^2 - \tilde{p}_1(\widehat{m})) + \frac{1}{\alpha} (\langle \mathbf{u}_{\widehat{m}}, \varepsilon \rangle_N^2 - \tilde{p}_2(\widehat{m})) + \frac{1}{\alpha} (\langle \mathbf{u}_m, \varepsilon \rangle_N^2 - \tilde{p}_2(m)) \right\}_+ \end{aligned}$$

using that  $m$  minimizes the function  $\left\| \mathbf{f} - \mathbf{P}_{m'} \mathbf{f} \right\|_N^2 + \text{pen}(m')$  and (5.15).

For any  $x > 0$ ,

$$\begin{aligned} \mathbb{P} \left( (1 - \alpha) \mathcal{H}(\mathbf{f}) \geq \frac{x \delta_m^2}{N} \right) &\leq \mathbb{P} \left( \exists m' \in \mathcal{M} : (2 - \alpha) (G_{m'}^2 - \tilde{p}_1(m')) \geq \frac{x \delta_{m'}^2}{3N} \right) \\ &\quad + \mathbb{P} \left( \exists m' \in \mathcal{M} : \frac{1}{\alpha} (\langle \mathbf{u}_{m'}, \varepsilon \rangle_N^2 - \tilde{p}_2(m')) \geq \frac{x \delta_{m'}^2}{3N} \right) \\ &\leq \sum_{m' \in \mathcal{M}} \mathbb{P} \left( (2 - \alpha) (\left\| \mathbf{P}_{m'} \varepsilon \right\|_N^2 - \tilde{p}_1(m')) \geq \frac{x \delta_{m'}^2}{3N} \right) \\ &\quad + \sum_{m' \in \mathcal{M}} \mathbb{P} \left( \frac{1}{\alpha} (\langle \mathbf{u}_{m'}, \varepsilon \rangle_N^2 - \tilde{p}_2(m')) \geq \frac{x \delta_{m'}^2}{3N} \right) \\ &:= \sum_{m' \in \mathcal{M}} P_{1,m'}(x) + \sum_{m' \in \mathcal{M}} P_{2,m'}(x). \end{aligned} \quad (5.16)$$

We first bound  $P_{2,m'}(x)$ . Let  $t$  be some positive number,

$$\mathbb{P} (|\langle \mathbf{u}_{m'}, \varepsilon \rangle_N| \geq t) \leq t^{-p} \mathbb{E} (|\langle \mathbf{u}_{m'}, \varepsilon \rangle_N|^p). \quad (5.17)$$

Since  $\langle \mathbf{u}_{m'}, \varepsilon \rangle_N = \frac{1}{N} \sum_{i=1}^N \langle \mathbf{u}_{im'}, \varepsilon_i \rangle$  with  $\varepsilon_i$  i.i.d. and with zero mean, then by Rosenthal's inequality we know that for some constant  $c(p)$  that depends on  $p$  only

$$\begin{aligned} c^{-1}(p) N^p \mathbb{E} |\langle \mathbf{u}_{m'}, \varepsilon \rangle_N|^p &\leq \sum_{i=1}^N \mathbb{E} |\langle \mathbf{u}_{im'}, \varepsilon_i \rangle|^p + \left( \sum_{i=1}^N \mathbb{E} (\langle \mathbf{u}_{im'}, \varepsilon_i \rangle^2) \right)^{\frac{p}{2}} \\ &\leq \sum_{i=1}^N \mathbb{E} \|\mathbf{u}_{im'}\|^p \|\varepsilon_i\|^p + \left( \sum_{i=1}^N \mathbb{E} \|\mathbf{u}_{im'}\|^2 \|\varepsilon_i\|^2 \right)^{\frac{p}{2}} \\ &= \mathbb{E} \|\varepsilon_1\|^p \sum_{i=1}^N \|\mathbf{u}_{im'}\|^p + (\mathbb{E} \|\varepsilon_1\|^2)^{\frac{p}{2}} \left( \sum_{i=1}^N \|\mathbf{u}_{im'}\|^2 \right)^{\frac{p}{2}}. \end{aligned} \quad (5.18)$$

Since  $p \geq 2$ ,  $(\mathbb{E} \|\varepsilon_1\|^2)^{\frac{1}{2}} \leq (\mathbb{E} \|\varepsilon_1\|^p)^{\frac{1}{p}}$  and

$$(\mathbb{E} \|\varepsilon_1\|^2)^{\frac{p}{2}} \leq \mathbb{E} \|\varepsilon_1\|^p. \quad (5.19)$$

Using also that by definition  $\|\mathbf{u}_{m'}\|_N^2 = \frac{1}{N} \sum_{i=1}^N \|\mathbf{u}_{im'}\|^2 = 1$ , then  $\frac{\|\mathbf{u}_{im'}\|^2}{N} \leq 1$  and therefore  $\frac{\|\mathbf{u}_{im'}\|}{N^{\frac{1}{2}}} \leq 1$ . Thus

$$\sum_{i=1}^N \|\mathbf{u}_{im'}\|^p = N^{\frac{p}{2}} \sum_{i=1}^N \left( \frac{\|\mathbf{u}_{im'}\|}{N^{\frac{1}{2}}} \right)^p \leq N^{\frac{p}{2}} \sum_{i=1}^N \left( \frac{\|\mathbf{u}_{im'}\|}{N^{\frac{1}{2}}} \right)^2 = N^{\frac{p}{2}} \|\mathbf{u}_{m'}\|_N^2 = N^{\frac{p}{2}}. \quad (5.20)$$

We deduce from (5.18), (5.19) and (5.20) that

$$c^{-1}(p) N^p \mathbb{E} |\langle \mathbf{u}_{m'}, \varepsilon \rangle_N|^p \leq \mathbb{E} \|\varepsilon_1\|^p N^{\frac{p}{2}} + \mathbb{E} \|\varepsilon_1\|^p N^{\frac{p}{2}}.$$

Then for some constant  $c'(p)$  that only depends on  $p$

$$\mathbb{E} |\langle \mathbf{u}_{m'}, \varepsilon \rangle_N|^p \leq c'(p) \mathbb{E} \|\varepsilon_1\|^p N^{-\frac{p}{2}}.$$

By this last inequality and (5.17) we get that

$$\mathbb{P} (|\langle \mathbf{u}_{m'}, \varepsilon \rangle_N| \geq t) \leq c'(p) \mathbb{E} \|\varepsilon_1\|^p N^{-\frac{p}{2}} t^{-p}. \quad (5.21)$$

Let  $v$  be some positive number depending on  $\eta$  only to be chosen later. We take  $t$  such that  $Nt^2 = \min(v, \frac{\alpha}{3}) (L_{m'} D_{m'} + x) \delta_{m'}^2$  and set  $N\tilde{p}_2(m') = v L_{m'} D_{m'} \delta_{m'}^2$ . We get

$$\begin{aligned} P_{2,m'}(x) &= \mathbb{P} \left( \frac{1}{\alpha} (\langle \mathbf{u}_{m'}, \varepsilon \rangle_N^2 - \tilde{p}_2(m')) \geq \frac{x \delta_{m'}^2}{3N} \right) \\ &= \mathbb{P} \left( N \langle \mathbf{u}_{m'}, \varepsilon \rangle_N^2 \geq N\tilde{p}_2(m') + \alpha \frac{\delta_{m'}^2}{3} x \right) \\ &= \mathbb{P} \left( N \langle \mathbf{u}_{m'}, \varepsilon \rangle_N^2 \geq v L_{m'} D_{m'} \delta_{m'}^2 + \alpha \frac{\delta_{m'}^2}{3} x \right) \\ &\leq \mathbb{P} \left( |\langle \mathbf{u}_{m'}, \varepsilon \rangle_N| \geq N^{-\frac{1}{2}} \sqrt{\min(v, \frac{\alpha}{3})} \sqrt{(L_{m'} D_{m'} + x) \delta_{m'}^2} \right) \\ &\leq c'(p) \mathbb{E} \|\varepsilon_1\|^p N^{-\frac{p}{2}} \frac{N^{\frac{p}{2}}}{(\min(v, \frac{\alpha}{3}))^{\frac{p}{2}} (L_{m'} D_{m'} + x)^{\frac{p}{2}} \delta_{m'}^p} \\ &= c''(p, \eta) \frac{\mathbb{E} \|\varepsilon_1\|^p}{\delta_m^p} \frac{1}{(L_{m'} D_{m'} + x)^{\frac{p}{2}}}. \end{aligned} \quad (5.22)$$

The last inequality holds using (5.21).

We now bound  $P_{1,m'}(x)$  for those  $m' \in \mathcal{M}$  such that  $D_{m'} \geq 1$ . By using our version of Corollary 5.1 in Baraud with  $\tilde{A} = \mathbf{P}_{m'}$ ,  $\text{Tr}(\tilde{A}) = D_{m'}$  and  $\rho(\tilde{A}) = 1$ , we obtain from (4.3) that for any positive  $x_{m'}$

$$\mathbb{P}\left(N \|\mathbf{P}_{m'}\varepsilon\|_N^2 \geq \delta_{m'}^2 D_{m'} + 2\delta_{m'}^2 \sqrt{D_{m'} x_{m'}} + \delta_{m'}^2 D_{m'} x_{m'}\right) \leq C(p) \frac{\mathbb{E} \|\varepsilon_1\|^p}{\delta_{m'}^p} D_{m'} x_{m'}^{-\frac{p}{2}}. \quad (5.23)$$

Since for any  $\beta > 0$ ,  $2\sqrt{D_{m'} x_{m'}} \leq \beta D_{m'} + \beta^{-1} x_{m'}$  then (5.23) imply that

$$\mathbb{P}\left(N \|\mathbf{P}_{m'}\varepsilon\|_N^2 \geq (1 + \beta) D_{m'} \delta_{m'}^2 + (1 + \beta^{-1}) x_{m'} \delta_{m'}^2\right) \leq C(p) \frac{\mathbb{E} \|\varepsilon_1\|^p}{\delta_{m'}^p} D_{m'} x_{m'}^{-\frac{p}{2}}. \quad (5.24)$$

Now for some number  $\beta$  depending on  $\eta$  only to be chosen later, we take  $x_{m'} = (1 + \beta^{-1}) \min\left(v, \frac{(2-\alpha)^{-1}}{3}\right) (L_{m'} D_{m'} + x)$  and  $N\tilde{p}_1(m') = v L_{m'} D_{m'} \delta_{m'}^2 + (1 + \beta) D_{m'} \delta_{m'}^2$ . By (5.24) this gives

$$\begin{aligned} P_{1,m'}(x) &= \mathbb{P}\left(\|\mathbf{P}_{m'}\varepsilon\|_N^2 - \tilde{p}_1(m') \geq \frac{(2-\alpha)^{-1} x \delta_{m'}^2}{3N}\right) \\ &= \mathbb{P}\left(N \|\mathbf{P}_{m'}\varepsilon\|_N^2 \geq v L_{m'} D_{m'} \delta_{m'}^2 + (1 + \beta) D_{m'} \delta_{m'}^2 + \frac{(2-\alpha)^{-1}}{3} x \delta_{m'}^2\right) \\ &\leq \mathbb{P}\left(N \|\mathbf{P}_{m'}\varepsilon\|_N^2 \geq (1 + \beta) D_{m'} \delta_{m'}^2 + (1 + \beta^{-1}) x_{m'} \delta_{m'}^2\right) \\ &\leq c(p) \frac{\mathbb{E} \|\varepsilon_1\|^p}{\delta_{m'}^p} D_{m'} x_{m'}^{-\frac{p}{2}} \leq c'(p, \eta) \frac{\mathbb{E} \|\varepsilon_1\|^p}{\delta_{m'}^p} \frac{D_{m'}}{(L_{m'} D_{m'} + x)^{\frac{p}{2}}}. \end{aligned} \quad (5.25)$$

Gathering (5.22), (5.25) and (5.16) we get that

$$\begin{aligned} \mathbb{P}\left(\mathcal{H}(\mathbf{f}) \geq \frac{x \delta_{m'}^2}{N(1-\alpha)}\right) &\leq \sum_{m' \in \mathcal{M}} P_{1,m'}(x) + \sum_{m' \in \mathcal{M}} P_{2,m'}(x) \\ &\leq \sum_{m' \in \mathcal{M}} c'(p, \eta) \frac{\mathbb{E} \|\varepsilon_1\|^p}{\delta_{m'}^p} \frac{D_{m'}}{(L_{m'} D_{m'} + x)^{\frac{p}{2}}} \\ &\quad + \sum_{m' \in \mathcal{M}} c''(p, \eta) \frac{\mathbb{E} \|\varepsilon_1\|^p}{\delta_{m'}^p} \frac{1}{(L_{m'} D_{m'} + x)^{\frac{p}{2}}}. \end{aligned}$$

Since  $\frac{1}{(1-\alpha)} = (1 + 2\eta^{-1})$ , then (5.3) holds:

$$\begin{aligned} \mathbb{P}\left(\mathcal{H}(\mathbf{f}) \geq (1 + 2\eta^{-1}) \frac{x \delta_{m'}^2}{N}\right) &\leq \sum_{m' \in \mathcal{M}} \frac{\mathbb{E} \|\varepsilon_1\|^p}{\delta_{m'}^p (L_{m'} D_{m'} + x)^{\frac{p}{2}}} \max(D_{m'}, 1) (c'(p, \eta) + c''(p, \eta)) \\ &= c(p, \eta) \frac{\mathbb{E} \|\varepsilon_1\|^p}{\delta_{m'}^p} \sum_{m' \in \mathcal{M}} \frac{D_{m'} \vee 1}{(L_{m'} D_{m'} + x)^{\frac{p}{2}}}. \end{aligned}$$

It remains to choose  $\beta$  and  $\delta$  for (5.14) to hold (we recall that  $\alpha = \frac{2}{2+\eta}$ ). This is the case if  $(2 - \alpha)(1 + \beta) = 1 + \eta$  and  $(2 - \alpha + \alpha^{-1})\delta = 1$ , therefore we take  $\beta = \frac{\eta}{2}$  and  $\delta = \left[1 + \frac{\eta}{2} + 2\frac{(1+\eta)}{(2+\eta)}\right]^{-1}$ .  $\square$

### 5.3 Proof of the concentration inequality

Proof of Proposition (4.3)

*Proof.* Denote by  $\tau^2$  the following expression:

$$\tau^2 := \mathbb{E} \|\mathbf{P}_m \varepsilon\|^2 = \mathbb{E} (\varepsilon^T \mathbf{P}_m \varepsilon) = \text{Tr} (\mathbf{P}_m (I_N \otimes \Phi)).$$

Then we have that

$$\begin{aligned} \tau^2 &= \text{Tr} (\mathbf{P}_m (I_N \otimes \Phi) \mathbf{P}_m) \leq \lambda_{\max} (I_N \otimes \Phi) \text{Tr} (\mathbf{P}_m^2) = \lambda_{\max} (I_N \otimes \Phi) \text{Tr} (\mathbf{P}_m) \\ &= \lambda_{\max} (\Phi) \text{Tr} (\mathbf{P}_m) = \lambda_{\max} (\Phi) D_m. \end{aligned}$$

We have that  $\eta^2(\varepsilon) := \varepsilon^T \tilde{A} \varepsilon$ , where  $\tilde{A} = A^T A$ . Then

$$\begin{aligned} \eta^2(\varepsilon) &= \|A\varepsilon\|^2 = \left[ \sup_{\|\mathbf{u}\| \leq 1} \langle A\varepsilon, \mathbf{u} \rangle \right]^2 = \left[ \sup_{\|\mathbf{u}\| \leq 1} \sum_{i=1}^{Nd} (A\varepsilon)_i \mathbf{u}_i \right]^2 \\ &= \left[ \sup_{\|\mathbf{u}\| \leq 1} \langle \varepsilon, A^T \mathbf{u} \rangle \right]^2 = \left[ \sup_{\|\mathbf{u}\| \leq 1} \sum_{i=1}^N \langle \varepsilon_i, (A^T \mathbf{u})_i \rangle \right]^2 \\ &= \left[ \sup_{\|\mathbf{u}\| \leq 1} \sum_{i=1}^N \langle \varepsilon_i, A_i^T \mathbf{u} \rangle \right]^2 = \left[ \sup_{\|\mathbf{u}\| \leq 1} \sum_{i=1}^N \sum_{j=1}^d \varepsilon_{ij} (A_i^T \mathbf{u})_j \right]^2 \end{aligned}$$

with  $A = (A_1 \mid \dots \mid A_N)$ , where  $A_i$  is a  $(Nd) \times d$  matrix.

Now take  $\mathcal{G} = \{g_{\mathbf{u}} : g_{\mathbf{u}}(\mathbf{x}) = \sum_{i=1}^N \langle \mathbf{x}_i, A_i^T \mathbf{u} \rangle = \sum_{i=1}^N \langle B_i \mathbf{x}, B_i A^T \mathbf{u} \rangle, \mathbf{u}, \mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_N)' \in \mathbb{R}^{(Nd)}, \|\mathbf{u}\| \leq 1\}$ .

Let  $M_i = [\mathbf{0}, \dots, \mathbf{0}, I_d, \mathbf{0}, \dots, \mathbf{0}]' \in \mathbb{R}^{(Nd) \times (Nd)}$ , where  $I_d$  is the  $i$ -th block of  $M_i$ ,  $B_i = [0, \dots, 0, I_d, 0, \dots, 0] \in \mathbb{R}^{(Nd) \times (Nd)}$ ,  $\varepsilon_i = B_i \varepsilon$  and  $M_i \varepsilon = [\mathbf{0}, \dots, \mathbf{0}, \varepsilon_i, \mathbf{0}, \dots, \mathbf{0}]'$ .

Then

$$\eta(\varepsilon) = \sup_{\|\mathbf{u}\| \leq 1} \sum_{i=1}^N g_{\mathbf{u}}(M_i \varepsilon).$$

Now take  $\mathbf{U}_i = M_i \varepsilon$ ,  $\varepsilon \in \mathbb{R}^{(Nd)}$ . Then for each positive number  $t$  and  $p > 0$

$$\begin{aligned} \mathbb{P}(\eta(\varepsilon) \geq \mathbb{E}(\eta(\varepsilon)) + t) &\leq \mathbb{P}(|\eta(\varepsilon) - \mathbb{E}(\eta(\varepsilon))| > t) \\ &\leq t^{-p} \mathbb{E}(|\eta(\varepsilon) - \mathbb{E}(\eta(\varepsilon))|^p) \text{ by Markov inequality} \\ &\leq c(p) t^{-p} \left\{ \mathbb{E} \left( \max_{i=1, \dots, N} \sup_{\|\mathbf{u}\| \leq 1} |\langle \varepsilon_i, A_i^T \mathbf{u} \rangle|^p \right) + \left[ \mathbb{E} \left( \sup_{\|\mathbf{u}\| \leq 1} \sum_{i=1}^N (\langle \varepsilon_i, A_i^T \mathbf{u} \rangle)^2 \right)^{p/2} \right]^2 \right\} \\ &= c(p) t^{-p} (\mathbb{E}_1 + \mathbb{E}_2^{p/2}). \end{aligned} \tag{5.26}$$

We start by bounding  $\mathbb{E}_1$ . For all  $\mathbf{u}$  such that  $\|\mathbf{u}\| \leq 1$  and  $i \in \{1, \dots, N\}$ ,

$$\|A_i^T \mathbf{u}\|^2 \leq \|A^T \mathbf{u}\|^2 \leq \rho^2(A),$$

where  $\rho(M) = \sup_{x \neq 0} \frac{\|Mx\|}{\|x\|}$  for all matrix  $M$ . For  $p \geq 2$  we have that  $\|A_i \mathbf{u}\|^p \leq \rho^{p-2}(A) \|A_i \mathbf{u}\|^2$ ,

then

$$|\langle \varepsilon_i, A_i^T \mathbf{u} \rangle|^p \leq [\|\varepsilon_i\| \|A_i^T \mathbf{u}\|]^p \leq \rho^{p-2}(A) \|\varepsilon_i\|^p \|A_i^T \mathbf{u}\|^p.$$

Therefore

$$\mathbb{E}_1 \leq \rho^{p-2}(A) \mathbb{E} \left( \sup_{\|\mathbf{u}\|=1} \sum_{i=1}^N \|\varepsilon_i\|^p \|A_i^T \mathbf{u}\|^2 \right).$$

Since  $\|\mathbf{u}\| \leq 1, \forall i = 1, \dots, N$

$$\begin{aligned} \|A_i^T \mathbf{u}\|^2 &= \mathbf{u}^T A_i A_i^T \mathbf{u} \leq \rho(A_i A_i^T) \\ &\leq \text{Tr}(A_i A_i^T), \end{aligned}$$

then

$$\sum_{i=1}^N \|A_i^T \mathbf{u}\|^2 \leq \sum_{i=1}^N \text{Tr}(A_i A_i^T) = \text{Tr} \left( \sum_{i=1}^N A_i A_i^T \right) = \text{Tr}(\tilde{A}).$$

Thus,

$$\mathbb{E}_1 \leq \rho^{p-2}(A) \text{Tr}(\tilde{A}) \mathbb{E}(\|\varepsilon_i\|^p). \quad (5.27)$$

We now bound  $\mathbb{E}_2$  via a truncation argument. Since for all  $\mathbf{u}$  such that  $\|\mathbf{u}\| \leq 1$  and  $i \in \{1, \dots, N\}$ ,  $\|A_i^T \mathbf{u}\|^2 \leq \rho^2(A)$ , for any positive number  $c$  to be specified later we have that

$$\begin{aligned} \mathbb{E}_2 &\leq \mathbb{E} \left( \sup_{\|\mathbf{u}\| \leq 1} \sum_{i=1}^N \|\varepsilon_i\|^2 \|A_i^T \mathbf{u}\|^2 1_{\{\|\varepsilon_i\| \leq c\}} \right) + \mathbb{E} \left( \sup_{\|\mathbf{u}\| \leq 1} \sum_{i=1}^N \|\varepsilon_i\|^2 \|A_i^T \mathbf{u}\|^2 1_{\{\|\varepsilon_i\| > c\}} \right) \\ &\leq \mathbb{E} \left( c^2 \sup_{\|\mathbf{u}\| \leq 1} \sum_{i=1}^N \|A_i^T \mathbf{u}\|^2 1_{\{\|\varepsilon_i\| \leq c\}} \right) + \mathbb{E} \left( \sup_{\|\mathbf{u}\| \leq 1} \sum_{i=1}^N \|\varepsilon_i\|^2 \|A_i^T \mathbf{u}\|^2 1_{\{\|\varepsilon_i\| > c\}} \right) \\ &\leq c^2 \rho^2(A) + c^{2-p} \mathbb{E} \left( \sup_{\|\mathbf{u}\| \leq 1} \sum_{i=1}^N \|A_i \mathbf{u}\|^2 \|\varepsilon_i\|^p \right) \\ &\leq c^2 \rho^2(A) + c^{2-p} \mathbb{E}(\|\varepsilon_i\|^p) \text{Tr}(\tilde{A}) \end{aligned} \quad (5.28)$$

using the bound obtained for  $\mathbb{E}_1$ . It remains to take  $c^p = \mathbb{E}(\|\varepsilon_i\|^p) \text{Tr}(\tilde{A}) / \rho^2(A)$  to get that:

$$\mathbb{E}_2 \leq c^2 \rho^2(A) + c^2 \rho^2(A) = 2c^2 \rho^2(A),$$

therefore

$$\mathbb{E}_2^{p/2} \leq 2^{p/2} c^p \rho^p(A), \quad (5.29)$$

which implies that

$$2^{-p/2} \mathbb{E}_2^{p/2} \leq \mathbb{E}(\|\varepsilon_1\|^p) \text{Tr}(\tilde{A}) \rho^{p-2}(A).$$

We straightforwardly derive from (5.26) that

$$\mathbb{P}(\eta^2(\varepsilon) \geq [\mathbb{E}(\eta(\varepsilon))]^2 + 2\mathbb{E}(\eta(\varepsilon))t + t^2) \leq c(p) t^{-p} (\mathbb{E}_1 + \mathbb{E}_2^{p/2}).$$

Since  $[\mathbb{E}(\eta(\varepsilon))]^2 \leq \mathbb{E}(\eta^2(\varepsilon))$ , (5.27) and (5.29) imply that

$$\begin{aligned} \mathbb{P}(\eta^2(\varepsilon) \geq \mathbb{E}(\eta^2(\varepsilon)) + 2\sqrt{\mathbb{E}(\eta^2(\varepsilon))}t + t^2) &\leq c(p) t^{-p} (\mathbb{E}_1 + \mathbb{E}_2^{p/2}) \\ &\leq c(p) t^{-p} (\rho^{p-2}(A) \text{Tr}(\tilde{A}) \mathbb{E}(\|\varepsilon_i\|^p) + 2^{p/2} \mathbb{E}(\|\varepsilon_1\|^p) \text{Tr}(\tilde{A}) \rho^{p-2}(A)) \\ &\leq c'(p) t^{-p} \rho^{p-2}(A) \text{Tr}(\tilde{A}) \mathbb{E}(\|\varepsilon_i\|^p), \end{aligned} \quad (5.30)$$

for all  $t > 0$ . Moreover

$$\begin{aligned}\mathbb{E}(\eta^2(\varepsilon)) &= \mathbb{E}(\varepsilon^T \tilde{A} \varepsilon) = \mathbb{E}(\|A\varepsilon\|^2) = \mathbb{E}\left(\sum_{i=1}^N \|A_i \varepsilon_i\|^2\right) \\ &= \sum_{i=1}^N \mathbb{E}(\text{Tr} \varepsilon_i^T A_i^T A_i \varepsilon_i) = \sum_{i=1}^N \text{Tr} A_i^T A_i \mathbb{E}(\varepsilon_i \varepsilon_i^T) \\ &= \text{Tr}\left(\sum_{i=1}^N A_i^T A_i\right) \Phi.\end{aligned}$$

But it is better to use that

$$\begin{aligned}\mathbb{E}(\eta^2(\varepsilon)) &= \text{Tr}(\tilde{A} \varepsilon \varepsilon^T) = \text{Tr}(\tilde{A}(I_N \otimes \Phi)) = \text{Tr}(A^T A (I_N \otimes \Phi)) = \text{Tr}(A(I_N \otimes \Phi) A^T) \\ &\leq \lambda_{\max}(I_N \otimes \Phi) \text{Tr}(A A^T) = \lambda_{\max}(I_N \otimes \Phi) \text{Tr}(\tilde{A}) = \lambda_{\max}(Q) \text{Tr}(\tilde{A}),\end{aligned}\tag{5.31}$$

for  $Q = I_N \otimes \Phi$ .

Using (5.31), take  $t^2 = \rho(\tilde{A}) \lambda_{\max}(I_N \otimes \Phi) x > 0$  in (5.30) to get that

$$\begin{aligned}\mathbb{P}\left(\eta^2(\varepsilon) \geq \lambda_{\max}(Q) \text{Tr}(\tilde{A}) + 2\sqrt{\lambda_{\max}(Q) \text{Tr}(\tilde{A}) \rho(\tilde{A}) \lambda_{\max}(Q) x} + \rho(\tilde{A}) \lambda_{\max}(Q) x\right) \\ \leq c'(p) \rho^{-p/2}(\tilde{A}) \lambda_{\max}^{-p/2}(Q) x^{-p/2} \rho^{p-2}(A) \text{Tr}(\tilde{A}) \mathbb{E}(\|\varepsilon_i\|^p).\end{aligned}$$

Since  $\rho(\tilde{A}) = \rho^2(A)$  (with the Euclidean norm) the desired result follows:

$$\begin{aligned}\mathbb{P}\left(\eta^2(\varepsilon) \geq \lambda_{\max}(Q) \text{Tr}(\tilde{A}) + 2\lambda_{\max}(Q) \sqrt{\rho(\tilde{A}) \text{Tr}(\tilde{A}) x} + \lambda_{\max}(Q) \rho(\tilde{A}) x\right) \\ \leq c'(p) \frac{\mathbb{E} \|\varepsilon_1\|^p}{\left(\sqrt{\lambda_{\max}(Q)}\right)^p} \frac{\text{Tr}(\tilde{A})}{\rho(\tilde{A}) x^{p/2}}.\end{aligned}\tag{5.32}$$

□

## References

- [Adl90] Robert J. Adler. *An introduction to continuity, extrema, and related topics for general Gaussian processes*. Institute of Mathematical Statistics Lecture Notes—Monograph Series, 12. Institute of Mathematical Statistics, Hayward, CA, 1990.
- [Bar00] Y. Baraud. Model selection for regression on a fixed design. *Probab. Theory Related Fields*, 117(4):467–493, 2000.
- [Bar02] Yannick Baraud. Model selection for regression on a random design. *ESAIM Probab. Statist.*, 6:127–146 (electronic), 2002.

- [BJG95] R. Biscay, J. C. Jimenez, and A. Gonzalez. Smooth approximation of nonnegative definite kernels. In *Approximation and optimization in the Caribbean, II (Havana, 1993)*, volume 8 of *Approx. Optim.*, pages 114–128. Lang, Frankfurt am Main, 1995.
- [BR97] Diaz-Frances E. Biscay, R. J. and L. M Rodriguez. Cross-validation of covariance structures using the frobenius matrix distance as a discrepancy function. *Journal of Statistical Computation and Simulation*, 1997.
- [Com01] Fabienne Comte. Adaptive estimation of the spectrum of a stationary Gaussian sequence. *Bernoulli*, 7(2):267–298, 2001.
- [Cre93] Noel A. C. Cressie. *Statistics for spatial data*. Wiley Series in Probability and Mathematical Statistics: Applied Probability and Statistics. John Wiley & Sons Inc., New York, 1993. Revised reprint of the 1991 edition, A Wiley-Interscience Publication.
- [EHN96] H. Engl, M. Hanke, and A. Neubauer. *Regularization of inverse problems*, volume 375 of *Mathematics and its Applications*. Kluwer Academic Publishers Group, Dordrecht, 1996.
- [ETA03] Perrin-O. Elogne, S. N. and C. Thomas-Agnan. Non parametric estimation of smooth stationary covariance functions by interpolation methods. *Phd*, 2003.
- [Gen08] Xavier Gendre. Simultaneous estimation of the mean and the variance in heteroscedastic Gaussian regression. *Electron. J. Stat.*, 2:1345–1372, 2008.
- [Jou77] A. G. Journel. Kriging in terms of projections. *J. Internat. Assoc. Mathematical Geol.*, 9(6):563–586, 1977.
- [KvR05] Tõnu Kollo and Dietrich von Rosen. *Advanced multivariate statistics with matrices*, volume 579 of *Mathematics and Its Applications (New York)*. Springer, Dordrecht, 2005.
- [LL08] J-M. Loubes and C. Ludena. Adaptive complexity regularization for inverse problems. *Electronic Journal Of Statistics*, 2:661–677, 2008.
- [LRZ08] Elizaveta Levina, Adam Rothman, and Ji Zhu. Sparse estimation of large covariance matrices via a nested Lasso penalty. *Ann. Appl. Stat.*, 2(1):245–263, 2008.
- [Lüt96] H. Lütkepohl. *Handbook of matrices*. John Wiley & Sons Ltd., Chichester, 1996.
- [MP08] Nychka D. W. Matsuo, T. and D. Paul. Nonstationary covariance modeling for incomplete data: smoothed monte-carlo approach. *preprint*, 2008.
- [RMSE08] N. Raj Rao, James A. Mingo, Roland Speicher, and Alan Edelman. Statistical eigen-inference from large Wishart matrices. *Ann. Statist.*, 36(6):2850–2885, 2008.
- [RS05] J. O. Ramsey and Silverman. *Functional Data Analysis*. Springer: NY, 2005.



- [SS05] Juliane Schäfer and Korbinian Strimmer. A shrinkage approach to large-scale covariance matrix estimation and implications for functional genomics. *Stat. Appl. Genet. Mol. Biol.*, 4:Art. 32, 28 pp. (electronic), 2005.
- [Ste99] Michael L. Stein. *Interpolation of spatial data. Some theory for kriging.* Springer Series in Statistics. New York, NY: Springer. xvii, 247 p., 1999.

J. BIGOT & J-M. LOUBES

Equipe de probabilités et statistique,  
Institut de Mathématique de Toulouse,  
UMR5219, Université de Toulouse,  
31000 Toulouse FRANCE

R. Biscay & L. Muñiz

Instituto de Cibernética, Matemática y Física,  
Departamento de Matemáticas,  
Universidad Central de la Habana,  
Ciudad Havana CUBA