



HAL
open science

PAC-Bayesian risk bounds for group-analysis sparse regression by exponential weighting

Tung Duy Luu, Jalal Fadili, Christophe Chesneau

► **To cite this version:**

Tung Duy Luu, Jalal Fadili, Christophe Chesneau. PAC-Bayesian risk bounds for group-analysis sparse regression by exponential weighting. 2016. hal-01367742v1

HAL Id: hal-01367742

<https://hal.science/hal-01367742v1>

Preprint submitted on 16 Sep 2016 (v1), last revised 22 May 2018 (v4)

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.

PAC-Bayesian risk bounds for group-analysis sparse regression by exponential weighting

Tung Duy Luu*

Jalal Fadili*

Christophe Chesneau†

Abstract

In this paper, we consider a high-dimensional non-parametric regression model with fixed design and i.i.d. random errors. We propose a powerful estimator by exponential weighted aggregation (EWA) with a group-analysis sparsity promoting prior on the weights. We prove that our estimator satisfies a sharp group-analysis sparse oracle inequality with a small remainder term ensuring its good theoretical performances. We also propose a forward-backward proximal Langevin Monte-Carlo algorithm to sample from the target distribution (which is not smooth nor log-concave) and derive its guarantees. In turn, this allows us to implement our estimator and validate it on some numerical experiments.

Key words. High-dimensional regression, sparse learning, exponential weighted aggregation, group-analysis sparsity, group-analysis sparse oracle inequality, frame, forward-backward Langevin Monte-Carlo algorithm

AMS subject classifications. 62G07 62G20

1 Introduction

1.1 Problem statement

Let us briefly present our statistical context. Assume that the given data (x_i, Y_i) , $i = 1, \dots, n$, is generated according to the non-parametric regression model

$$Y_i = f(x_i) + \xi_i, \quad i \in \{1, \dots, n\}, \quad (1.1)$$

where x_1, \dots, x_n are deterministic in an arbitrary set \mathcal{X} , $f : \mathcal{X} \rightarrow \mathbb{R}$ is the unknown regression function and (ξ_1, \dots, ξ_n) is a random error. (1.1) is equivalently written in vector form $\mathbf{Y} = \mathbf{f} + \boldsymbol{\xi}$. Assume that there exists a dictionary $\mathcal{H} = \{f_j : \mathcal{X} \rightarrow \mathbb{R}, j \in \{1, \dots, M\}\}$ such that f is well approximated by a linear combination of elements in \mathcal{H} , i.e., there exists $\tilde{\boldsymbol{\theta}} = (\tilde{\theta}_1, \dots, \tilde{\theta}_M)^T \in \mathbb{R}^M$ such that $f_{\tilde{\boldsymbol{\theta}}} = \sum_{j=1}^M \tilde{\theta}_j f_j$ is a suitable approximation of f . The f_j 's are known and may be either fixed atoms in a basis or pre-estimators. In the context of high dimension, the cardinality of \mathcal{H} is much larger than the sample size (i.e., $M \gg n$). Thus, in such a setting, the classical least-squares to estimate $\tilde{\boldsymbol{\theta}}$ is obviously not applicable.

The idea of aggregating elements in a dictionary has been introduced in machine learning to combine different techniques (see [40, 55]) with some procedures such as bagging [8], boosting [28, 52] and random forests [1, 5–7, 9, 29]. In the recent years, there has been a flurry of research on the use of the concept of

*Normandie Univ, ENSICAEN, UNICAEN, CNRS, GREYC, France, Email: {duy-tung.luu, Jalal.Fadili}@ensicaen.fr.

†Normandie Univ, UNICAEN, CNRS, LMNO, France, Email: christophe.chesneau@unicaen.fr.

sparsity in various areas including statistics and machine learning in high dimension. The idea is that even if the cardinality of \mathcal{H} is very large, the number of effective elements in the dictionary is much smaller than the sample size. Namely, the number of non-zero components of $\tilde{\theta}$ is assumed much smaller than n . This makes it possible to build an estimate $f_{\tilde{\theta}}$ with good provable performance guarantees under appropriate conditions on the dictionary and noise.

1.2 Overview of previous work

1.2.1 Oracle inequalities

This type of guarantees dates back, for instance, to the work [22–24] on orthogonal wavelet thresholding estimators. Oracle inequalities, which are at the heart of our work, quantify the quality of an estimator compared to the best possible one that could only be given with an oracle. These inequalities are well adapted under the sparsity scenario. Formally, let $g : \mathcal{X} \rightarrow \mathbb{R}$, and denote

$$\|g\|_n = \sqrt{\frac{1}{n} \sum_{j=1}^n g^2(x_j)}.$$

The performance of an estimator \hat{f} is measured by its averaged squared error, i.e.,

$$R(\hat{f}) = \|\hat{f} - f\|_n^2.$$

We aim to find an estimator \hat{f} that mimics as much as possible the best model of aggregation in a given class Θ (in the probabilistic sense). This idea is expressed in the following type of inequalities:

$$\mathbb{E} \left[\|\hat{f} - f\|_n^2 \right] \leq C \inf_{\theta \in \Theta} \left[\|f_{\theta} - f\|_n^2 + \Delta_{n,M}(\theta) \right], \quad (1.2)$$

where $C \geq 1$ and the remainder term $\Delta_{n,M}(\theta)$ depends on the performance of the estimator, the complexity of θ , the dimension M and the sample size n . Such type of inequality is called balanced oracle inequality. Under the sparsity scenario, the complexity of θ is characterized by the number of its non-zero components, in which case inequalities of type (1.2) are called sparse oracle inequalities (SOI).

An estimator with good oracle properties would correspond to C is close to 1 (ideally, $C = 1$, in which case the inequality is said “sharp”), $\Delta_{n,M}(\theta)$ is small even if $n \ll M$ and decreases rapidly to 0 as $n \rightarrow +\infty$. Besides, the choice of Θ is crucial: on the one hand, a non suitable choice can lead a large bias term in (1.2). On the other hand, if Θ is too complex, the remainder term becomes large. Then, a suitable choice for Θ must achieve a good bias-complexity trade-off.

In the literature, there are mainly two approaches to provide aggregated estimators in high dimension under the sparsity assumption: Penalization and Exponential Weighted Aggregation (EWA). Given $\mathbf{Y} = \mathbf{y}$, The penalization approach considers the minimization problem

$$\min_{\theta \in \Theta} \|\mathbf{y} - f_{\theta}\|_n^2 + \text{pen}(\theta),$$

where $\text{pen} : \mathbb{R}^M \rightarrow \mathbb{R}^+$ is a sparsity promoting penalty function, see e.g. [10]. Our work focuses on EWA approach that we briefly describe now.

1.2.2 Exponential Weighted Aggregation (EWA)

Let (Λ, \mathcal{A}) a space equipped with a σ -algebra and

$$\mathcal{F}_\Lambda = \{f_\lambda : \mathcal{X} \rightarrow \mathbb{R} : \lambda \in \Lambda\}$$

be a given collection (\mathcal{F}_Λ is called dictionary of aggregation) where $\lambda \rightarrow f_\lambda(x)$ is measurable $\forall x \in \mathcal{X}$. The functions f_λ may be deterministic or random. The aggregators depend on the nature of f_λ if the latter is random. Otherwise, the aggregators are defined via the probability measure

$$\mu_n(d\lambda) = \frac{\exp\left(-n\|\mathbf{Y} - f_\lambda\|_n^2/\beta\right)\pi(d\lambda)}{\int_\Lambda \exp\left(-n\|\mathbf{Y} - f_\omega\|_n^2/\beta\right)\pi(d\omega)}, \quad (1.3)$$

where $\beta > 0$ called temperature parameter and π called prior which is a probability measure on Λ . Then, we define the aggregate by

$$\hat{f}_n(x) = \int_\Lambda f_\lambda(x)\mu_n(d\lambda). \quad (1.4)$$

This idea was initially proposed in [34, 40, 55] with a uniform prior on a finite set Λ .

In the literature, the fact that the elements in the dictionary may be either deterministic or random splits the EWA approach into two corresponding cases. The deterministic case is considered in [15, 16, 18, 20]. These papers proposed several PAC-Bayesian type of oracle inequalities under different assumptions. Especially, the assumptions in [20] depend only on the noise and turns out to be fulfilled for a large class of noise. This serves to construct, for a suitable prior and dictionary, a SOI with a remainder term of order $O(\|\boldsymbol{\theta}\|_0 \log(M)/n)$, which scales linearly with the sparsity level and increases in M only logarithmically.

The random case is tackled in [44]. The initial idea is to obtain two independent samples from the initial sample by randomization or sample splitting (see [38, 47, 57]). The first sample is used to construct the pre-estimators, and the aggregation is performed on the second sample conditionally on the first one. However this idea does not work when the observations are not i.i.d. Several authors have proposed exponentially aggregating linear pre-estimators without splitting, and with discrete priors on the weights. Typical cases of linear pre-estimators are orthogonal projectors on all possible linear subspaces that are in the model set (e.g. in the sparsity context, linear subspaces spanned by the standard basis restricted to supports of increasing size). This was introduced in [39]. More recent work such as [17] generalizes the idea where the pre-estimators are affine and the priors are continuous.

1.2.3 Generalization of sparsity assumption

Analysis sparsity Let $P \geq M$ and $\mathbf{D} : \mathbb{R}^M \rightarrow \mathbb{R}^P$ be a linear analysis operator. The analysis sparsity assumption means that $\|\mathbf{D}\boldsymbol{\theta}\|_0 \ll n$. A typical example is total variation [51] where the operator \mathbf{D} corresponds to “finite differences” (i.e., $(\mathbf{D}\boldsymbol{\theta})_1 = \theta_1$, $(\mathbf{D}\boldsymbol{\theta})_j = \theta_j - \theta_{j-1}$, $\forall j \geq 2$). Another example is the fused Lasso [54] where \mathbf{D} is a positive combination of the identity and finite differences.

Group sparsity Group sparsity corresponds to saying that the aggregator $\boldsymbol{\theta}$ is block sparse, see Section 4.1 for formal details. Group sparsity is at the heart of the group Lasso and related methods [32, 36, 42, 43, 46, 58]. In the EWA context, the group sparsity prior is considered in [48] as an application of the aggregation of orthogonal projectors.

1.3 Contributions

Our main contributions are summarized as follows:

- We propose an EWA estimator, with a deterministic dictionary, under a group-analysis sparsity prior (see Section 1.2.3). More precisely, we assume that \mathbf{D} is the analysis operator corresponding to a frame of \mathbb{R}^M , and thus is not necessarily invertible unlike previous work. In addition, our prior class (see (5.2)) is much more general than previously proposed ones [20] which recovered as very special cases. It also allows more flexibility to enhance the performance of EWA. This prior class is parameterized through a function $g : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ that satisfies mild conditions (see Assumption 5.2).
- We establish a group-analysis SOI where the remainder term depends on the number of active groups in $\mathbf{D}\boldsymbol{\theta}$ and on the function g (see Theorem 6.1)
- For an appropriate choice of g which is well-adapted to the group-analysis sparsity scenario, we exhibit a group-analysis SOI where, as expected, the remainder term scales as $O\left(\|\mathbf{D}\boldsymbol{\theta}\|_{0,\mathcal{G}} \log(L)/n\right)$, where $\|\mathbf{D}\boldsymbol{\theta}\|_{0,\mathcal{G}}$ is the number of active groups in $\mathbf{D}\boldsymbol{\theta}$, and L is the total number of groups (see Corollary 6.2 and Remark 6.2). This rate coincides with the classical one $O(\|\boldsymbol{\theta}\|_0 \log(M)/n)$ under the sparsity scenario, i.e. $\mathbf{D} = \mathbf{I}_M$ and $L = M$.
- We also propose a forward-backward proximal Langevin Monte-Carlo (LMC) algorithm to sample from the target distribution (which is not smooth nor log-concave), and establish several of its properties in a general setting. In turn, this allows us to efficiently implement our EWA estimator with the proposed prior. We validate this algorithm on some numerical examples.

1.4 Paper organization

Necessary notations and some preliminaries are first introduced in Section 2. Section 3 reminds the PAC-Bayesian type oracle inequalities proposed in [20] which are a classical starting point in literature for EWA in the deterministic case. Section 4 consists in specifying our framework, i.e., group-analysis sparsity scenario where the analysis operator \mathbf{D} constitute a frame. In Section 5, we describe our EWA procedure after specifying the aggregation dictionary and our prior family. In Section 6, we establish our main results, namely group-analysis SOI. Section 7 is devoted to the forward-backward proximal LMC algorithm that implements EWA, and the numerical experiments on several numerical settings are described in Section 8. The proofs of all results are collected in Section 10.

2 Preliminaries

Throughout the paper, we will use the following notations.

For each function $h : \mathcal{X} \rightarrow \mathbb{R}$, we define $\mathbf{h} = (h(x_1), \dots, h(x_n))^T$ its n -sample vector form and

$$\|h\|_n = \left(\frac{1}{n} \sum_{j=1}^n h^2(x_j) \right)^{1/2}$$

its empirical norm. Moreover, for each function $h : \mathbb{R}^M \rightarrow \mathbb{R}$, we denote

$$\|h\|_\infty = \sup_{\mathbf{x} \in \mathbb{R}^M} |h(\mathbf{x})|.$$

For all $\mathbf{v} \in \mathbb{R}^d$, $d \geq 1$, we also define the ℓ_p function

$$\|\mathbf{v}\|_p = \left(\sum_{j=1}^d |v_j|^p \right)^{1/p}, \quad p > 0,$$

with the usual adaptation $\|\mathbf{v}\|_\infty = \max_{j \in \{1, \dots, d\}} |v_j|$. It is a norm for $p \geq 1$, and a quasi-norm for $p \in]0, 1[$. $\|\mathbf{v}\|_0$ is the ℓ_0 pseudo-norm which counts the number of non-zero elements in \mathbf{v} .

Let $\mathbf{A} \in \mathbb{R}^{p \times q}$, we set $\boldsymbol{\sigma}(\mathbf{A}) = (\sigma_1(\mathbf{A}), \dots, \sigma_q(\mathbf{A}))^T \in \mathbb{R}^q$ be the vector of singular values of \mathbf{A} in non-increasing order. Note that, when \mathbf{A} is symmetric semi-definite positive, $\boldsymbol{\sigma}(\mathbf{A})$ is also the ordered vector of positive eigenvalues of \mathbf{A} . The spectral norm of \mathbf{A} is

$$\|\mathbf{A}\| = \sqrt{\sigma_1(\mathbf{A}^T \mathbf{A})}.$$

For a set Ω , I_Ω is its characteristic function, i.e. 1 if the argument is in Ω and 0 otherwise.

Recall the Gamma function $\Gamma :]0, +\infty[\rightarrow]0, +\infty[$

$$\Gamma(t) = \int_0^\infty x^{t-1} \exp(-x) dx.$$

The two following lemmas contain useful formula used throughout the paper.

Lemma 2.1 ([30, 3.251.11]). *Let $p, \gamma, \nu, \eta > 0$. If $\gamma/\nu < \eta + 1$ we have*

$$\int_0^\infty \frac{x^{\gamma-1}}{(p+x^\nu)^{\eta+1}} dx = \frac{1}{\nu p^{\eta+1-\gamma/\nu}} \frac{\Gamma(\gamma/\nu)\Gamma(1+\eta-\gamma/\nu)}{\Gamma(1+\eta)}, \quad (2.1)$$

otherwise this integral is not definite.

Lemma 2.2 (Cartesian to spherical coordinates [27]). *Let $d \geq 1$ and a mapping $h : \mathbb{R}^+ \rightarrow \mathbb{R}$ such that $\mathbf{u} \rightarrow h(\|\mathbf{u}\|_2)$ is measurable in \mathbb{R}^d . We obtain*

$$\int_{\mathbb{R}^d} h(\|\mathbf{u}\|_2) d\mathbf{u} = C_d \int_0^\infty x^{d-1} h(x) dx, \quad (2.2)$$

where $C_d = 2\pi^{d/2}/\Gamma(d/2)$ is the surface area of a d -dimensional ball of radius 1.

3 PAC-Bayesian type oracle inequalities

This section recalls a PAC-Bayesian type oracle inequality which holds for the EWA procedure of type (1.3)-(1.4) with any deterministic aggregation dictionary, any prior and a large class of noises. Such type of oracle inequalities was introduced in [20] for i.i.d. noise. In the present paper, we adapt it to the non i.i.d. case. Indeed, let us start with the two following assumptions.

Assumption 3.1 (Standard assumption). *The noise vector $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n)^T$ has zero mean.*

Assumption 3.2 (Main assumption). *For any $\gamma > 0$ small enough, there exist a probability space and two random variables $\boldsymbol{\xi}'$ and ζ defined on this probability space satisfying the three following points:*

1. ξ' has the same distribution as ξ .
2. $\xi' + \zeta$ has the same distribution as $(1 + \gamma)\xi'$ and the conditional expectation satisfies $\mathbb{E}[\zeta|\xi'] = 0$.
3. There exist $t_0 \in (0, \infty]$ and a bounded Borel function $v : \mathbb{R}^M \rightarrow \mathbb{R}^+$ such that

$$\limsup_{\gamma \rightarrow 0} \sup_{(\mathbf{t}, \mathbf{a}) \in \mathbb{R}^M \times \mathbb{R}^M : (\|\mathbf{t}\|_2, \mathbf{a}) \in [-t_0, t_0] \times \text{supp}(\xi')}$$

$$\frac{\log \mathbb{E}[\exp(\mathbf{t}^T \zeta) | \xi' = \mathbf{a}]}{\|\mathbf{t}\|_2^2 \gamma v(\mathbf{a})} \leq 1.$$

where $\text{supp}(\xi')$ is the support of the distribution of ξ' .

Assumption 3.2 is based on [20, Assumption N], and can be shown to be fulfilled for a large class of noise such as Gaussian, Laplace and any bounded symmetric noise.

Besides, let $H \in]0, +\infty]$ such that

$$\sup_{(\boldsymbol{\lambda}, \boldsymbol{\lambda}') \in \Lambda^2} \|\mathbf{f}_{\boldsymbol{\lambda}} - \mathbf{f}_{\boldsymbol{\lambda}'}\|_2 \leq H. \quad (3.1)$$

Note that (3.1) is always satisfied since H is allowed to be infinite. However, for the sake of sharpness in our theoretical results, we wish to choose H as small as possible. We are now ready to state the PAC-Bayesian type oracle inequalities.

Theorem 3.1. *Let Assumptions 3.1 and 3.2 be satisfied with some function v and let (3.1) holds. Then for any prior π , any probability measure p on Λ and any $\beta \geq \max(4\|v\|_\infty, 2H/t_0)$ or $\beta \geq 4\|v\|_\infty$ when $H = \infty$, $t_0 = \infty$, we have*

$$\mathbb{E} \left[\|\widehat{f}_n - f\|_n^2 \right] \leq \int_{\Lambda} \|f - f_{\boldsymbol{\lambda}}\|_n^2 p(d\boldsymbol{\lambda}) + \frac{\beta \text{KL}(p, \pi)}{n}, \quad (3.2)$$

where \widehat{f}_n is the aggregate defined in (1.4) and $\text{KL}(p, \pi) = \int_{\Lambda} \log(p(d\boldsymbol{\lambda})/\pi(d\boldsymbol{\lambda}))p(d\boldsymbol{\lambda})$ is the Kullback-Leibler divergence.

The proof of Theorem 3.1 is a mild adaptation of the original one in [20, Section 2], where we used directly Assumption 3.2-3 in the vector ζ instead of splitting it into $\zeta_{i, i \in \{1, \dots, n\}}$ (that are no longer i.i.d.).

Related work The work of [16] proposed three types of oracle inequalities which are similar to (3.2) under different assumptions. The first type (see [16, Theorem 1]) holds under a restrictive condition on the noise. The second (see [16, Theorem 2]) involves conditions depending on the noise and also on the dictionary. The last (see [16, Theorem 4]) works for all symmetric noises without conditions on the dictionary. However, an additional term appears in the remainder term which has a low rate for some types of noise. Therefore, Theorem 3.1 (with Assumption 3.2) is a good trade-off between these types of oracle inequalities.

Moreover, there exist some related forms of (3.2) in different frameworks. For example, when $\xi_i \sim \mathcal{N}(0, \sigma_i^2)$, $i = 1, \dots, n$, the following aggregate was proposed in [17]:

$$\widehat{f} = \int_{\Lambda} \widehat{f}_{\boldsymbol{\lambda}} p(d\boldsymbol{\lambda}), \quad p(d\boldsymbol{\lambda}) = \frac{\exp\left(-\frac{n}{\beta} \widehat{r}_{\boldsymbol{\lambda}}\right) \pi(d\boldsymbol{\lambda})}{\int_{\Lambda} \exp\left(-\frac{n}{\beta} \widehat{r}_{\boldsymbol{\omega}}\right) \pi(d\boldsymbol{\omega})},$$

where $\widehat{f}_{\boldsymbol{\lambda}}, \boldsymbol{\lambda} \in \Lambda$ are affine estimators satisfying some conditions imposed in [17, Theorem 1] which yield the definition of $\widehat{r}_{\boldsymbol{\lambda}}, \boldsymbol{\lambda} \in \Lambda$. This aggregate satisfies oracle inequalities defined therein which are the counterparts of (3.2) for the aggregation of estimators. In addition, in the case of random design (i.e., x_1, \dots, x_n are random and i.i.d.), the works in [19] constructed a mirror averaging aggregate to obtain a generalized type of oracle inequalities where the performance is measured by any loss instead of the averaged square loss.

4 Group-analysis sparsity

4.1 Group sparsity

We now describe formally what is intended by group-analysis sparsity, which measures group sparsity of the image of a vector with an analysis linear sparsifying transform.

Let $P \geq M$. We partition the index set $\{1, \dots, P\}$ into L non-overlapping groups/blocks of indices $\{\mathcal{G}_l\}_{1 \leq l \leq L}$ such that

$$\mathcal{G} = \bigcup_{l=1}^L \mathcal{G}_l = \{1, \dots, P\} \quad \text{and} \quad \mathcal{G}_l \cap \mathcal{G}_k = \emptyset, \quad \forall l \neq k.$$

For the sake of simplicity, and without loss of generality, the groups are assumed to have the same size $\text{Card } \mathcal{G}_l = G \geq 1$ and the total number of blocks L is supposed to be an integer. A vector $\mathbf{v} \in \mathbb{R}^P$ can be divided into L vectors $\mathbf{v}_{\mathcal{G}_l} \in \mathbb{R}^G$ which are the restrictions of \mathbf{v} to the coordinates indexed by \mathcal{G}_l . We define

$$\|\mathbf{v}\|_{p,\mathcal{G}} = \left(\sum_{l=1}^L \|\mathbf{v}_{\mathcal{G}_l}\|_2^p \right)^{1/p}, \quad p > 0,$$

which is a norm for $p \geq 1$, with $\|\mathbf{v}\|_{\infty,\mathcal{G}} = \max_{l \in \{1, \dots, L\}} \|\mathbf{v}_{\mathcal{G}_l}\|_2$. It is a quasi-norm for $p \in]0, 1[$. $\|\mathbf{v}\|_{0,\mathcal{G}}$ counts the number of active (i.e. non-zero) groups in \mathbf{v} . With these notations, the group-analysis sparsity assumption is formalized as follows.

Assumption 4.1 (Group-analysis sparsity assumption). *Let $\mathbf{D} : \mathbb{R}^M \rightarrow \mathbb{R}^P$ be a linear analysis operator. There exists $\tilde{\boldsymbol{\theta}} \in \mathbb{R}^M$ such that $f_{\tilde{\boldsymbol{\theta}}}$ is close to f and $\|\mathbf{D}\tilde{\boldsymbol{\theta}}\|_{0,\mathcal{G}} \ll n$.*

In plain words, Assumption 4.1 says that the number of active groups of $\mathbf{D}\tilde{\boldsymbol{\theta}}$ is much smaller than the sample size. Note that this is a strict notion of group-analysis sparsity, and a weaker one could be also considered where most $(\mathbf{D}\tilde{\boldsymbol{\theta}})_{\mathcal{G}_l}$ are nearly zero.

4.2 Frame assumption

Let $\mathbf{D} : \mathbb{R}^M \rightarrow \mathbb{R}^P$ be a linear operator seen as a matrix in $\mathbb{R}^{P \times M}$. We assume the following throughout this paper.

Assumption 4.2 (Frame assumption). *\mathbf{D} corresponds to the analysis operator of a finite frame of \mathbb{R}^M , i.e. there exist two constants ν and μ with $\nu \geq \mu > 0$, called frame bounds, such that the generalized Parseval relation is satisfied*

$$\mu \|\boldsymbol{\theta}\|_2^2 \leq \|\mathbf{D}\boldsymbol{\theta}\|_2^2 \leq \nu \|\boldsymbol{\theta}\|_2^2, \quad \forall \boldsymbol{\theta} \in \mathbb{R}^M.$$

By the Courant-Fischer theorem, Assumption 4.2 is equivalent to the fact that μ (resp. ν) is a lower (resp. upper) bound of the eigenvalues of $\mathbf{D}^T \mathbf{D}$. Moreover, since $\mu > 0$, we have that

$$\mathbf{D}^T \mathbf{D} \text{ is bijective and } \mathbf{D} \text{ is injective.}$$

The frame is said tight when $\mu = \nu$. Typical examples of (tight) frames that have been used in statistics are translation invariant wavelets [14], ridgelets [12] and curvelets [11].

Let $\widetilde{\mathbf{D}}^T \in \mathbb{R}^{M \times P}$ be the matrix whose columns form the canonical dual frame associated to \mathbf{D}^T . We know that

$$\widetilde{\mathbf{D}}^T \mathbf{D} = \mathbf{I}_M \quad (4.1)$$

and

$$\frac{1}{\mu} \geq \sigma_1 \left(\widetilde{\mathbf{D}}^T \widetilde{\mathbf{D}}^{T^T} \right) \geq \dots \geq \sigma_M \left(\widetilde{\mathbf{D}}^T \widetilde{\mathbf{D}}^{T^T} \right) \geq \frac{1}{\nu}. \quad (4.2)$$

Note that we focus on the canonical dual frame for the sake of simplicity. In fact, our exposition in the rest of the paper remains unchanged if any other dual frame is used instead of the canonical one.

Under the frame Assumption 4.2, the following lemma provides an efficient change of variables formula, which will be a key tool in the proof of our general group-analysis SOI (see Theorem 6.1).

Lemma 4.1. *Let $\Theta \subseteq \mathbb{R}^M$ be a measurable set. For any $\mathbf{D} \in \mathbb{R}^{P \times M}$ satisfying Assumption 4.2, let $u : \mathbb{R}^P \rightarrow \mathbb{R}$ such that the mapping $\boldsymbol{\theta} \mapsto u(\mathbf{D}\boldsymbol{\theta})$ is measurable on Θ . We have*

$$\int_{\Theta} u(\mathbf{D}\boldsymbol{\theta}) d\boldsymbol{\theta} = \frac{1}{\sqrt{\det(\mathbf{D}^T \mathbf{D})}} \int_{\mathbf{D}\Theta} u(\mathbf{v}) d\mathbf{v} \quad (4.3)$$

provided either u is non-negative valued or the integral on the left converges.

Though quite natural, proving Lemma 4.1 rigorously requires nontrivial arguments from geometric measure theory (see Section 10.1.1 for details).

5 EWA estimator

To design an aggregation by exponential weighting, two ingredients are essential: the aggregation dictionary and the prior which promotes group-analysis sparsity. We specify them below.

5.1 Choice of dictionary

Let $\mathbf{X} \in \mathbb{R}^{n \times M}$ where $\mathbf{X}_{i,j} = f_j(x_i)$, $f_j \in \mathcal{H}$. We impose the following standard normalization assumption on \mathbf{X} .

Assumption 5.1. *\mathbf{X} is normalized such that all the diagonal entries of $\mathbf{X}^T \mathbf{X} / n$ are 1.*

Now, let us introduce our dictionary of aggregation:

$$\mathcal{F}_{\Theta} = \left\{ f_{\boldsymbol{\theta}} = \ell \left(\sum_{j=1}^M \theta_j f_j \right) : \boldsymbol{\theta} \in \Theta \right\}, \quad (5.1)$$

where $\Theta = \left\{ \boldsymbol{\theta} \in \mathbb{R}^M : \|\mathbf{D}\boldsymbol{\theta}\|_{a,\mathcal{G}}^a \leq R \right\}$, $a \in]0, 1]$, $R \in]0, +\infty]$ and $\ell : \mathbb{R} \rightarrow \mathbb{R}$ is twice continuously differentiable and known depending on the regression problem (for example: $\ell(x) = e^x$ for the exponential regression, $\ell(x) = e^x / (e^x + 1)$ for the logistic regression and $\ell(x) = x$ for the linear regression). This dictionary of aggregation is similar to the one proposed in [18–20]. However, the set of index is modified to adapt the group-analysis sparsity and the exponent a is varied in $]0, 1]$ instead of a fixed $a = 1$. The bound H in (3.1) for \mathcal{F}_{Θ} in (5.1) is established in the following result.

Proposition 5.1. Let $\mathcal{F}_\Theta = \{f_\theta : \theta \in \Theta\}$ defined in (5.1) with some $R > 0$, $a \in]0, 1]$ and $\ell : \mathbb{R} \rightarrow \mathbb{R}$ twice continuously differentiable. Let Assumption 4.2 holds for some $\mu > 0$. We get that

$$\sup_{(\theta, \theta') \in \Theta^2} \|f_\theta - f_{\theta'}\|_2 \leq 2 \max_{x \in \mathcal{B}} \|L(x)\|_2,$$

where $\mathcal{B} = \{x \in \mathbb{R}^n : \|x\|_2 \leq \|X\| R^{1/a} / \sqrt{\mu}\}$ and $L : x \in \mathbb{R}^n \rightarrow (\ell(x_1), \dots, \ell(x_n))$.

From Proposition 5.1, one can choose $H = 2 \max_{x \in \mathcal{B}} \|L(x)\|_2$.

5.2 Choice of prior

We choose a general prior of the form

$$\pi(d\theta) = \frac{1}{C_{\alpha, g, R}} \prod_{l=1}^L \exp\left(-\alpha^a \|[D\theta]_{\mathcal{G}_l}\|_2^a\right) g\left(\|[D\theta]_{\mathcal{G}_l}\|_2\right) I_\Theta(\theta) d\theta, \quad (5.2)$$

where $\alpha \geq 0$ and g satisfies the following requirements:

Assumption 5.2.

1. $g : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ is a bounded function such that $g \not\equiv 0$, $\theta \mapsto g(\|[D\theta]_{\mathcal{G}_l}\|_2)$ is measurable on \mathbb{R}^M , $\forall l \in \{1, \dots, L\}$.
2. The integrability condition:

$$\int_{\mathbb{R}^M} \prod_{l=1}^L g(\|[Du]_{\mathcal{G}_l}\|_2) du < \infty. \quad (5.3)$$

3. The moment condition:

$$\int_{\mathbb{R}^M} \|[Du]_{\mathcal{G}_l}\|_2^{2a} \prod_{k=1}^L g(\|[Du]_{\mathcal{G}_k}\|_2) du < \infty, \quad \forall l \in \{1, \dots, L\}. \quad (5.4)$$

4. There exist $\lambda \geq 0$ and a function $h : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ such that $\forall (\mathbf{t}, \mathbf{t}^*) \in \mathbb{R}^G \times \mathbb{R}^G$,

$$\frac{g(\|\mathbf{t} - \mathbf{t}^*\|_2)}{g(\|\mathbf{t}\|_2)} \leq h(\|\mathbf{t}^*\|_2)^\lambda.$$

Assumptions 5.2-2 and 5.2-3 lead the following remark that is a core part for the construction of the general group-analysis SOI in Theorem 6.1.

Remark 5.1. Let $G \geq 1$ and $D \in \mathbb{R}^{P \times M}$ satisfying Assumption 4.2. For any function g satisfying Assumptions 5.2-1, 5.2-2 and 5.2-3, and any $a \in]0, 1]$, there exists $K_{a, g} \in]0, \infty[$ such that, $\forall l \in \{1, \dots, L\}$,

$$\frac{\int_{\mathbb{R}^M} \|[Du]_{\mathcal{G}_l}\|_2^{2a} \prod_{k=1}^L g(\|[Du]_{\mathcal{G}_k}\|_2) du}{\int_{\mathbb{R}^M} \prod_{k=1}^L g(\|[Dv]_{\mathcal{G}_k}\|_2) dv} \leq K_{a, g}. \quad (5.5)$$

Proof. Remark 5.1 is briefly proved as follows. From Assumption 5.2-3 and the fact that $g \not\equiv 0$, one can show that (5.5) holds for $a = 1$. Moreover, since $g \not\equiv 0$ and g satisfies Assumption 5.2-2, we have

$$\mathbf{u} \rightarrow \frac{\prod_{l=1}^L g(\|[\mathbf{D}\mathbf{u}]_{\mathcal{G}_l}\|_2) d\mathbf{u}}{\int_{\mathbb{R}^M} \prod_{k=1}^L g(\|[\mathbf{D}\mathbf{v}]_{\mathcal{G}_k}\|_2) d\mathbf{v}}$$

is a probability measure. Therefore, (5.5) holds for any a in $]0, 1]$ by Hölder's inequality. \square

At first glance, Assumptions 5.2-2 and 5.2-3 may seem cumbersome. However the following lemma gives a simple condition on g that implies them.

Lemma 5.1. *Let $G \geq 1$ and $\mathbf{D} \in \mathbb{R}^{P \times M}$ satisfying Assumption 4.2. Suppose that g satisfies Assumption 5.2-1 and*

$$\int_0^\infty z^{G+1} g(z) dz < \infty. \quad (5.6)$$

Then Assumptions 5.2-2 and 5.2-3 are in force.

In the two following remarks, we consider the case where \mathbf{D} is invertible. Remark 5.2 provides a simple and explicit form for $K_{a,g}$ and Remark 5.3 shows that condition (5.6) is necessary for g to obey Assumption 5.2.

Remark 5.2. *Let $G \geq 1$ and $\mathbf{D} \in \mathbb{R}^{M \times M}$ be invertible. For any function g satisfying Assumptions 5.2-1, 5.2-2 and 5.2-3, and for any $a \in]0, 1]$, one can choose $K_{a,g}$ in (5.5) as*

$$K_{a,g} = \frac{\int_0^\infty x^{G-1+2a} g(x) dx}{\int_0^\infty z^{G-1} g(z) dz}. \quad (5.7)$$

Proof. The proof follows by combining Lemmas 4.1 and 2.2, i.e.,

$$\begin{aligned} \frac{\int_{\mathbb{R}^M} \|[\mathbf{D}\mathbf{u}]_{\mathcal{G}_i}\|_2^{2a} \prod_{k=1}^L g(\|[\mathbf{D}\mathbf{u}]_{\mathcal{G}_k}\|_2) d\mathbf{u}}{\int_{\mathbb{R}^M} \prod_{k=1}^L g(\|[\mathbf{D}\mathbf{v}]_{\mathcal{G}_k}\|_2) d\mathbf{v}} &= \frac{\int_{\mathbb{R}^G} \|\mathbf{u}\|_2^{2a} g(\|\mathbf{u}\|_2) d\mathbf{u}}{\int_{\mathbb{R}^G} g(\|\mathbf{v}\|_2) d\mathbf{v}} \\ &= \frac{\int_0^\infty x^{G-1+2a} g(x) dx}{\int_0^\infty y^{G-1} g(y) dy}. \end{aligned}$$

\square

Remark 5.3. *When \mathbf{D} is invertible, if g does not satisfy (5.6) then g cannot fulfill Assumption 5.2-3. Consequently, Assumption 5.2-3 and condition (5.6) are equivalent in the invertible case.*

Proof. By Lemmas 4.1 and 2.2, we get

$$\int_{\mathbb{R}^M} \|[\mathbf{D}\mathbf{u}]_{\mathcal{G}_i}\|_2^2 \prod_{k=1}^L g(\|[\mathbf{D}\mathbf{u}]_{\mathcal{G}_k}\|_2) d\mathbf{u} = \frac{C_G^L \int_0^\infty z^{G+1} g(z) dz}{\sqrt{\det(\mathbf{D}^T \mathbf{D})}} \left(\int_0^\infty w^{G-1} g(w) dw \right)^{L-1}.$$

\square

Let us now discuss some choices of g . The goal is to find a prior leading an oracle inequality with a small remainder term while promoting group-analysis sparsity. Namely, g should be peaked around 0 with heavy tails.

Example 5.1. Consider $g : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ defined by

$$g(x) = \frac{1}{(\tau^2 + x^2)^2}, \quad \tau > 0.$$

This choice of g yields a prior that specializes to the one in [20] for the individual sparsity scenario, i.e. with $\mathbf{D} = \mathbf{I}_M$, $G = 1$ and $a = 1$.

Example 5.2. Consider $g : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ defined by

$$g(x) = \frac{1}{(\tau^b + x^b)^c},$$

where $\tau > 0$, $b \in]0, 1]$ and $c > (2+G)/b$. The choice of c guarantees the validity of Assumption 5.2. Thanks to the parameters b and c , this choice of g offers more flexibility than the one of the previous example. This allows for example to optimize the performance of EWA by tuning these parameters for the particular dataset at hand.

6 Group-analysis sparse oracle inequality

Once a suitable dictionary and prior are chosen according to the above, the EWA is performed via (1.3)-(1.4). Our goal now is to provide a theoretical guarantee for the aggregates by constructing a group-analysis SOI. First of all, based on PAC-Bayesian type oracle inequalities in Section 3, we establish our first main result: a group-analysis SOI for the dictionary (5.1) and the prior (5.2) with a function g obeying Assumptions 5.2-1 to 4.

Theorem 6.1 (General group-analysis sparse oracle inequality). *Let $G \geq 1$ and \mathbf{D} satisfying Assumption 4.2 with $\mu > 0$. Let Assumptions 3.1 and 3.2 be satisfied with some function v , (3.1) holds and $\beta \geq \max(4\|v\|_\infty, 2H/t_0)$. For some $a \in]0, 1]$, take the dictionary (5.1) and the prior (5.2) with g satisfying Assumptions 5.2-1 to 4. Let $K_{a,g}$, as defined in (5.5), and assume that $R > 3\sqrt{K_{a,g}L}$. Then the following group-analysis SOI holds,*

$$\begin{aligned} \mathbb{E} \left[\|\widehat{f}_n - f\|_n^2 \right] \leq & \inf_{\boldsymbol{\theta} \in \mathbb{R}^M : \|\mathbf{D}\boldsymbol{\theta}\|_{a,g}^a \leq R - 3\sqrt{K_{a,g}L}} \left[\|f_{\boldsymbol{\theta}} - f\|_n^2 \right. \\ & + \frac{\beta}{n} \left(1 + 3\sqrt{K_{a,g}L}\alpha^a + \alpha^a \|\mathbf{D}\boldsymbol{\theta}\|_{a,g}^a \right) \\ & + \left. \frac{\lambda\beta}{n} \sum_{l=1}^L \log \left\{ h \left(\|\mathbf{D}\boldsymbol{\theta}\|_{G_l,2} \right) \right\} \right] \\ & + \frac{2K_{1,g}e^{3\sqrt{K_{a,g}L}\alpha^a} MC_{f,\ell}}{\mu}, \end{aligned} \tag{6.1}$$

where $C_{f,\ell} = \|\ell'\|_\infty^2 + \|\ell''\|_\infty (\|\ell\|_\infty + \|f\|_\infty)$.

See Section 10.2.1 for the proof.

In this group-analysis SOI, the EWA estimator \widehat{f}_n mimics the best aggregate $f_{\boldsymbol{\theta}}$ for all possible weights belonging to $\left\{ \boldsymbol{\theta} \in \mathbb{R}^M : \|\mathbf{D}\boldsymbol{\theta}\|_{a,g}^a \leq R - 3\sqrt{K_{a,g}L} \right\}$. The remainder term of the inequality (6.1) depends

on the choice of g via the function h and $K_{a,g}$. Theorem 6.1 has several consequences that we report now, for instance with the choices of g in Example 5.1 and 5.2.

We first consider the prior (5.2) in Example 5.1, under the individual sparsity scenario ($\mathbf{D} = \mathbf{I}_M, G = 1$) and the choice $a = 1$ (i.e., $\Theta = \{\boldsymbol{\theta} \in \mathbb{R}^M : \|\boldsymbol{\theta}\|_1 \leq R\}$). This is the setting considered in [20]. We obtain the following SOI as a corollary of our main result.

Corollary 6.1. *Let $\mathbf{D} = \mathbf{I}_M$ and fix $G = 1$. Suppose that Assumptions 3.1 and 3.2 hold with some function v , (3.1) holds and $\beta \geq \max(4\|v\|_\infty, 2H/t_0)$. Fix $a = 1$, take the dictionary (5.1) and the prior (5.2) with g defined in Example 5.1 and $\alpha \leq 1/(3M\tau)$. Assume that $R > 3M\tau$. Then,*

$$\begin{aligned} \mathbb{E} \left[\|\widehat{f}_n - f\|_n^2 \right] \leq & \inf_{\boldsymbol{\theta} \in \mathbb{R}^M : \|\boldsymbol{\theta}\|_1 \leq R - 3M\tau} \left[\|f\boldsymbol{\theta} - f\|_n^2 + \frac{2\beta}{n} \left(1 + \frac{1}{6M\tau} \|\boldsymbol{\theta}\|_1 \right) \right. \\ & \left. + \frac{4\beta}{n} \sum_{j=1}^M \log \left\{ 1 + \frac{|\theta_j|}{\tau} \right\} \right] + 2\tau^2 C_{f,\ell} M e, \end{aligned} \quad (6.2)$$

where $C_{f,\ell} = \|\ell'\|_\infty^2 + \|\ell''\|_\infty (\|\ell\|_\infty + \|f\|_\infty)$.

SOI (6.2) is similar to the one in [20, Theorem 2]. By choosing $\tau^2 \sim 1/(Mn)$ and $R \sim M\tau$, its remainder term is of order $O(\|\boldsymbol{\theta}\|_0 \log(M)/n)$ which is the classical rate under the sparsity scenario. However, the following remark shows that this prior is not adapted in the group-analysis case for any group size strictly larger than 1.

Remark 6.1. *Suppose that $G \geq 2$, and let $\gamma = G + 2$, $\nu = 2$ and $\eta = 1$. We have $\gamma/\nu \geq \eta + 1$, and thus Lemma 2.1 yields*

$$\int_0^\infty \frac{x^{G+1}}{(\tau^2 + x^2)^2} dx \text{ is not definite.}$$

Consequently, condition (5.6) is not fulfilled with g defined in Example 5.1 when $G \geq 2$.

According to Remark 5.3, Remark 6.1 implies that Assumption 5.2 is not fulfilled for g in Example 5.1 when the group size $G \geq 2$ and \mathbf{D} invertible. Thus one cannot construct a group-analysis SOI from Theorem 6.1 to guarantee the quality of the corresponding estimator. Overcoming this limitation was yet another motivation behind the choice of g in Example 5.2, which turns out to work well under the group-analysis sparsity scenario. In a nutshell, an aggregate with g in Example 5.2 exhibits the group-analysis SOI defined in the following corollary with any $G \geq 1$, any $\mathbf{D} \in \mathbb{R}^{P \times M}$ satisfying Assumption 4.2 and any $a \in]0, 1]$.

Corollary 6.2. *Let $G \geq 1$ and \mathbf{D} satisfying Assumption 4.2 with $\mu > 0$. Let Assumptions 3.1 and 3.2 be satisfied with some function v , (3.1) holds and $\beta \geq \max(4\|v\|_\infty, 2H/t_0)$. Take the dictionary (5.1) and the prior (5.2) with $a \in]0, 1]$, $\alpha \geq 0$ and g defined in Example 5.2. We get that g satisfies Assumption 5.2. Then,*

let $K_{a,g}$ as defined in (5.5), and assume that $R > 3\sqrt{K_{a,g}}L$. Then the following group-analysis SOI holds,

$$\begin{aligned}
\mathbb{E} \left[\|\widehat{f}_n - f\|_n^2 \right] &\leq \inf_{\boldsymbol{\theta} \in \mathbb{R}^M: \|\mathbf{D}\boldsymbol{\theta}\|_{a,\mathcal{G}}^a \leq R-3\sqrt{K_{a,g}}L} \left[\|f\boldsymbol{\theta} - f\|_n^2 \right. \\
&+ \frac{\beta}{n} \left(1 + 3\sqrt{K_{a,g}}L\alpha^a + \alpha^a \|\mathbf{D}\boldsymbol{\theta}\|_{a,\mathcal{G}}^a \right) \\
&+ \left. \frac{c\beta}{n} \sum_{l=1}^L \log \left(1 + \left[\frac{\|[\mathbf{D}\boldsymbol{\theta}]_{\mathcal{G}_l}\|_2}{\tau} \right]^b \right) \right] \\
&+ \frac{2K_{1,g}e^{3\sqrt{K_{a,g}}L\alpha^a} MC_{f,\ell}}{\mu}, \tag{6.3}
\end{aligned}$$

where $C_{f,\ell} = \|\ell'\|_\infty^2 + \|\ell''\|_\infty (\|\ell\|_\infty + \|f\|_\infty)$.

To get an explicit control of the remainder term, it is instructive to have a closed-form of $K_{a,g}$. This can be done for instance when \mathbf{D} is invertible, see (5.7). The obtained group-analysis SOI is stated as follows.

Corollary 6.3. *Consider the same framework as Corollary 6.2 with \mathbf{D} invertible. For $a \in]0, 1]$, let $\widetilde{K}_{a,g} = \frac{\Gamma((2a+G)/b)\Gamma(c-(2a+G)/b)}{\Gamma(G/b)\Gamma(c-G/b)}$, and set $\alpha \leq 1/\left(3\tau^a\sqrt{\widetilde{K}_{a,g}}L\right)^{1/a}$. Then we obtain the following group-analysis SOI*

$$\begin{aligned}
\mathbb{E} \left[\|\widehat{f}_n - f\|_n^2 \right] &\leq \inf_{\boldsymbol{\theta} \in \mathbb{R}^M: \|\mathbf{D}\boldsymbol{\theta}\|_{a,\mathcal{G}}^a \leq R-3\tau^a\sqrt{\widetilde{K}_{a,g}}L} \left[\|f\boldsymbol{\theta} - f\|_n^2 \right. \\
&+ \frac{c\beta}{n} \sum_{l=1}^L \log \left(1 + \left[\frac{\|[\mathbf{D}\boldsymbol{\theta}]_{\mathcal{G}_l}\|_2}{\tau} \right]^b \right) \\
&+ \left. \frac{\beta}{n} \left(2 + \frac{\|\mathbf{D}\boldsymbol{\theta}\|_{a,\mathcal{G}}^a}{3L\tau^a\sqrt{\widetilde{K}_{a,g}}} \right) \right] + \frac{2C_{f,\ell}M\tau^2\widetilde{K}_{1,g}e}{\mu}, \tag{6.4}
\end{aligned}$$

where $C_{f,\ell} = \|\ell'\|_\infty^2 + \|\ell''\|_\infty (\|\ell\|_\infty + \|f\|_\infty)$.

For an efficient aggregate, the parameters should be chosen such that the remainder term is small and decreases rapidly to 0 as $n \rightarrow \infty$. A suitable choice is described in the following remark.

Remark 6.2. *Set $\tau^2 \sim 1/(Mn)$ and $R \sim L\tau^a$, we have*

$$\begin{aligned}
\frac{\beta}{n} \left(2 + \frac{\|\mathbf{D}\boldsymbol{\theta}\|_{a,\mathcal{G}}^a}{3L\tau^a\sqrt{\widetilde{K}_{a,g}}} \right) + \frac{2C_{f,\ell}M\tau^2\widetilde{K}_{1,g}e}{\mu} &\leq \frac{\beta}{n} \left(2 + \frac{R}{3L\tau^a\sqrt{\widetilde{K}_{a,g}}} \right) + \frac{2C_{f,\ell}M\tau^2\widetilde{K}_{1,g}e}{\mu} \\
&= O\left(\frac{1}{n}\right),
\end{aligned}$$

and

$$\begin{aligned} \frac{c\beta}{n} \sum_{l=1}^L \log \left(1 + \left[\frac{\|[\mathbf{D}\boldsymbol{\theta}]_{\mathcal{G}_l}\|_2}{\tau} \right]^b \right) &\leq \frac{c\beta}{n} \sum_{l=1}^L \log \left(1 + \left[\frac{R^{1/a}}{\tau} \right]^b \right) \\ &= O\left(\frac{\|\mathbf{D}\boldsymbol{\theta}\|_{0,\mathcal{G}} \log(L)}{n} \right). \end{aligned} \quad (6.5)$$

In (6.5), $\|\mathbf{D}\boldsymbol{\theta}\|_{0,\mathcal{G}}$ is small for instance when $\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}$, by the group-analysis sparsity assumption (R must be sufficiently large to cover $\tilde{\boldsymbol{\theta}}$).

From Remark 6.2, the aggregate with g in Example 5.2 performs well under the group-analysis sparsity scenario. Under the sparsity scenario, this rate becomes $O(\|\boldsymbol{\theta}\|_0 \log(M)/n)$ which is the same rate as the aggregate with g in Example 5.1.

7 Forward-Backward proximal LMC algorithm

The goal of this section is to implement our EWA estimator (with g defined in Example 5.2) via a novel forward-backward proximal Monte-Carlo algorithm based on the Langevin diffusion (coined FB-LMC).

Let us consider a linear regression problem where $\ell(x) = x$, and thus

$$\hat{f}_n = \mathbf{X}\hat{\boldsymbol{\theta}}_n,$$

where $\mathbf{X} \in \mathbb{R}^{n \times M}$ is the design matrix described in Section 5.1. We also consider the EWA estimator

$$\hat{\boldsymbol{\theta}}_n = \int_{\mathbb{R}^M} \boldsymbol{\theta} \mu_n(d\boldsymbol{\theta}),$$

with the prior (5.2), i.e.

$$\mu_n(d\boldsymbol{\theta}) \propto \exp\left(\frac{-\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2}{\beta} \right) \pi(d\boldsymbol{\theta}). \quad (7.1)$$

Computing $\hat{\boldsymbol{\theta}}_n$ corresponds to an integration problem which becomes very involved to solve analytically or even numerically in high-dimension. A classical alternative is to approximate it via a Markov chain Monte-Carlo (MCMC) method which consists in sampling from μ by constructing an appropriate Markov chain whose stationary distribution is μ , and to compute sample path averages based on the output of the Markov chain. The theory of MCMC methods is based on that of Markov chains on continuous state space. As in [20], we here use the Langevin diffusion process; see [49].

7.1 The Langevin diffusion

Continuous dynamics A Langevin diffusion \mathbf{L} in \mathbb{R}^M , $M \geq 1$ is a homogeneous Markov process defined by the stochastic differential equation (SDE)

$$d\mathbf{L}(t) = \frac{1}{2}\boldsymbol{\rho}(\mathbf{L}(t))dt + d\mathbf{W}(t), \quad t > 0, \quad \mathbf{L}(0) = \mathbf{l}_0, \quad (7.2)$$

where $\boldsymbol{\rho} = \nabla \log \nu$, ν is everywhere non-zero and suitably smooth target density function on \mathbb{R}^M , \mathbf{W} is a M -dimensional Brownian process and $\mathbf{l}_0 \in \mathbb{R}^M$ is the initial value. Under mild assumptions, the SDE

(7.2) has a unique strong solution and, moreover, $\mathbf{L}(t)$ has a stationary distribution with density precisely ν [49, Theorem 2.1]. $\mathbf{L}(t)$ is therefore interesting for sampling from ν . In particular, this opens the door to approximating integrals $\int_{\mathbb{R}^M} F(\boldsymbol{\theta})\nu(\boldsymbol{\theta})d\boldsymbol{\theta}$ by the average value of a Langevin diffusion, i.e., $\frac{1}{T} \int_0^T F(\mathbf{L}(t))dt$ for a large enough T . Under additional assumptions on ν and G in a proper functional class, the expected squared error of the approximation can be controlled [56].

Forward Euler discretization In practice, in simulating the diffusion sample path, we cannot follow exactly the dynamic defined by the SDE (7.2). Instead, we must discretize it. A popular discretization is given by the forward (Euler) scheme, which reads

$$\mathbf{L}_{k+1} = \mathbf{L}_k + \frac{\delta}{2}\boldsymbol{\rho}(\mathbf{L}_k) + \sqrt{\delta}\mathbf{Z}_k, \quad t > 0, \quad \mathbf{L}_0 = \mathbf{l}_0,$$

where $\delta > 0$ is a sufficiently small constant discretization step-size and $\{\mathbf{Z}_k\}_k$ are i.i.d. $\sim \mathcal{N}(0, \mathbf{I}_M)$. The average value $\frac{1}{T} \int_0^T \mathbf{L}(t)dt$ can then be naturally approximated via the Riemann sum

$$\frac{\delta}{T} \sum_{k=0}^{\lfloor T/\delta \rfloor} \mathbf{L}_k.$$

It is then natural to approximate $\widehat{\boldsymbol{\theta}}_n$ by applying this discretization strategy to the Langevin diffusion with μ in (7.1) as the target density. However, quantitative consistency guarantees of this discretization require ν (hence $\boldsymbol{\rho}$) to be sufficiently smooth, which limits their applicability in our context. To cope with this difficulty, several works have proposed to replace $\log \nu$ with a smoothed version (typically involving the Moreau-Yosida regularization/envelope, see Definition 7.2) [20, 25, 26, 45]. In [26, 45] for instance, the authors proposed proximal-type algorithms to sample from possibly non-smooth log-concave densities ν using the forward Euler discretization and the Moreau-Yosida regularization. In [45]¹, $-\log \nu$ is replaced with its Moreau envelope, while in [26], it is assumed that $-\log \nu = F + G$, F is convex Lipschitz continuously differentiable, and G is a proper closed convex function replaced by its Moreau envelope. In both these works, convexity plays a crucial role to get quantitative convergence guarantees and thus cannot be applied to our prior. Proximal steps within MCMC methods have been recently proposed for some simple (convex) signal processing problems [13], though without any guarantees.

7.2 Forward-Backward proximal LMC algorithm

Recall that in our context, μ_n in (1.3) with the prior (5.2) is not differentiable nor log-concave. To overcome these difficulties, we will exploit the structure of μ_n and involved arguments from variational analysis [50]. For ease of notation, in the following, we denote with the same symbol the measure and its density with respect to the Lebesgue measure. Let

$$V(\boldsymbol{\theta}) \stackrel{\text{def}}{=} -\log \mu_n(\boldsymbol{\theta}) = F(\boldsymbol{\theta}) + G(\boldsymbol{\theta}), \quad (7.3)$$

where

$$F(\boldsymbol{\theta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2/\beta \quad \text{and} \quad G(\boldsymbol{\theta}) = -\log \pi(\boldsymbol{\theta}).$$

¹The author however applied it to problems where $-\log \nu = F + G$. But the gradient of the Moreau envelope of a sum, which amounts to computing the proximity operator of $-\log \nu$ does not have an easily implementable expression even if those of F and G do.

It is important to emphasize that the results we derive in the rest of the section apply to a general class of potentials G that encompasses as a special case the one corresponding to the prior π (5.2). In Section 7.3, we show that these results indeed apply to (5.2).

7.2.1 A glimpse of variational analysis

Before proceeding, we need a bit of notations and definitions. A more comprehensive account on variational analysis in finite-dimensional Euclidean spaces can be found in [50].

Definition 7.1 (Subdifferential). *Given a point $\boldsymbol{\theta} \in \mathbb{R}^M$ where a function $V : \mathbb{R}^M \rightarrow \mathbb{R} \cup \{+\infty\}$ is finite, the subdifferential of V at $\boldsymbol{\theta}$ is defined as,*

$$\partial V(\boldsymbol{\theta}) = \{\mathbf{v} \in \mathbb{R}^M : \exists \boldsymbol{\theta}_k \rightarrow \boldsymbol{\theta}, V(\boldsymbol{\theta}_k) \rightarrow V(\boldsymbol{\theta}), \mathbf{v} \leftarrow \mathbf{v}_k \in \partial^F V(\boldsymbol{\theta}_k)\},$$

where the Fréchet subdifferential $\partial^F V(\boldsymbol{\theta})$ of V at $\boldsymbol{\theta}$, is the set of vectors \mathbf{v} such that

$$V(\mathbf{w}) \geq V(\boldsymbol{\theta}) + (\mathbf{v})^T (\mathbf{w} - \boldsymbol{\theta}) + o(\|\mathbf{w} - \boldsymbol{\theta}\|_2).$$

V is subdifferentially regular at $\boldsymbol{\theta}$ if and only if V is locally lower-semicontinuous (lsc) there with $\partial V(\boldsymbol{\theta}) = \partial^F V(\boldsymbol{\theta})$.

$\partial V(\boldsymbol{\theta})$ and $\partial^F V(\boldsymbol{\theta})$ are closed, with $\partial^F V(\boldsymbol{\theta})$ convex and $\partial^F V(\boldsymbol{\theta}) \subset \partial V(\boldsymbol{\theta})$ [50, Theorem 8.6]. A proper lsc convex function is subdifferentially regular.

A function V is proper if it is not identically $+\infty$ and $V(\boldsymbol{\theta}) > -\infty$ for all $\boldsymbol{\theta}$.

Definition 7.2 (Proximal mapping and Moreau envelope). *Let $\mathbf{M} \in \mathbb{R}^{M \times M}$ be symmetric positive definite. For a proper lsc function V and $\gamma > 0$, the proximal mapping and Moreau envelope in the metric \mathbf{M} are defined respectively by*

$$\begin{aligned} \text{prox}_{\gamma V}^{\mathbf{M}}(\boldsymbol{\theta}) &\stackrel{\text{def}}{=} \underset{\mathbf{w} \in \mathbb{R}^M}{\text{Argmin}} \frac{1}{2\gamma} \|\mathbf{w} - \boldsymbol{\theta}\|_{\mathbf{M}}^2 + V(\mathbf{w}) \\ \mathbf{M}, \gamma V(\boldsymbol{\theta}) &\stackrel{\text{def}}{=} \inf_{\mathbf{w} \in \mathbb{R}^M} \frac{1}{2\gamma} \|\mathbf{w} - \boldsymbol{\theta}\|_{\mathbf{M}}^2 + V(\mathbf{w}) \end{aligned}$$

$\text{prox}_{\gamma V}^{\mathbf{M}}$ here is a set-valued operator since the minimizer, if it exists, is not necessarily unique. When $\mathbf{M} = \mathbf{I}_M$, we simply write $\text{prox}_{\gamma V}$ and γV .

V is prox-bounded if there exists $\gamma > 0$ such that $\mathbf{M}, \gamma V(\boldsymbol{\theta}) > -\infty$ for some $\boldsymbol{\theta}$. The supremum of the set of all such γ is the threshold of prox-boundedness for V .

Definition 7.3 (Hypomonotone and monotone operators). *A set-valued operator $S : \mathbb{R}^M \rightrightarrows \mathbb{R}^M$ is hypomonotone of modulus $r > 0$ if*

$$(\boldsymbol{\theta}' - \boldsymbol{\theta})^T (\boldsymbol{\eta}' - \boldsymbol{\eta}) \geq -r \|\boldsymbol{\theta}' - \boldsymbol{\theta}\|_2^2, \quad \forall \boldsymbol{\theta}', \boldsymbol{\theta} \in \mathbb{R}^M \text{ and } \boldsymbol{\eta}' \in S(\boldsymbol{\theta}'), \boldsymbol{\eta} \in S(\boldsymbol{\theta}).$$

It is monotone if the inequality holds with $r = 0$.

7.2.2 Key properties

We start by collecting some important properties on the potential V in (7.3) that will be crucial on our way to design our algorithm. To avoid trivialities, from now on, we assume that $\text{Argmin}(V) \neq \emptyset$. We also state the following assumptions.

(H.1) $G : \mathbb{R}^M \rightarrow \mathbb{R} \cup \{+\infty\}$ is proper, lsc and bounded from below.

(H.2) $\gamma \in]0, \frac{\beta}{2\|\mathbf{X}\|^2}[$.

Lemma 7.1. *Assume that (H.1) and (H.2) hold. Define $M_\gamma \stackrel{\text{def}}{=} \mathbf{I}_M - (2\gamma/\beta)\mathbf{X}^T\mathbf{X}$, which is symmetric positive definite. Then $\text{prox}_{\gamma V}^{M_\gamma}(\boldsymbol{\theta})$ and $\text{prox}_{\gamma G}(\boldsymbol{\theta})$ are non-empty compact sets for any $\boldsymbol{\theta}$. Moreover,*

$$\text{prox}_{\gamma V}^{M_\gamma} = \text{prox}_{\gamma G} \circ (\mathbf{I}_M - \gamma \nabla F),$$

and

$$\boldsymbol{\theta} \in \text{Argmin}(V) \Rightarrow \boldsymbol{\theta} \in \text{prox}_{\gamma V}^{M_\gamma}(\boldsymbol{\theta}).$$

The careful reader may have recognized from this lemma the classical forward-backward operator associated to $F + G$. This lemma states that this operator coincides with the proximal mapping of γV in the metric M_γ . Moreover, fixed points of this proximal mapping include minimizers of V . They are not equal however in general, unless for instance V is convex.

Lemma 7.2. *Assume that (H.1) and (H.2) hold. Then $M_{\gamma, \gamma} V(\boldsymbol{\theta})$ is finite and depends continuously on $(\boldsymbol{\theta}, \gamma) \in \mathbb{R}^M \times]0, \frac{\beta}{2\|\mathbf{X}\|^2}[$, and $(M_{\gamma, \gamma} V(\boldsymbol{\theta}))_{\gamma \in]0, \frac{\beta}{2\|\mathbf{X}\|^2}[}$ is a decreasing net. More precisely,*

$$M_{\gamma, \gamma} V(\boldsymbol{\theta}) \nearrow V(\boldsymbol{\theta}) \text{ for all } \boldsymbol{\theta} \text{ as } \gamma \searrow 0.$$

In the following lemma, we need the additional assumption:

(H.3) $\text{prox}_{\gamma G}$ is single-valued on \mathbb{R}^M , for any γ satisfying (H.2).

Lemma 7.3. *Assume that (H.1), (H.2) and (H.3) hold. Then $\text{prox}_{\gamma V}^{M_\gamma}$ is single-valued, continuous in $(\boldsymbol{\theta}, \gamma) \in \mathbb{R}^M \times]0, \frac{\beta}{2\|\mathbf{X}\|^2}[$, and $M_{\gamma, \gamma} V \in C^1(\mathbb{R}^M)$ with gradient*

$$\nabla M_{\gamma, \gamma} V = \gamma^{-1} M_\gamma \left(\mathbf{I}_M - \text{prox}_{\gamma V}^{M_\gamma} \right).$$

In plain words, Lemma 7.1 and 7.3 tell us that the action of the operator $\text{prox}_{\gamma G} \circ (\mathbf{I}_M - \gamma \nabla F)$ is equivalent to a gradient descent on the Moreau envelope of V in the metric M_γ with step-size γ .

To afford even better properties of $M_{\gamma, \gamma} V$, we need to strengthen our assumptions on G and γ .

(H.4) ∂G is hypomonotone of modulus r .

(H.5) $\gamma \in]0, \frac{1}{r+2\|\mathbf{X}\|^2/\beta}[$.

It is obvious that (H.5) implies (H.2). The hypomonotonicity assumption is equivalent to saying that $G + \frac{r}{2} \|\cdot\|_2^2$ is convex, or that G has a supporting quadratic minorant at each point of its domain. Such functions are closely related to lower- C^2 and $1/r$ -proximal functions in variational analysis; see [50, Theorem 10.33 and Example 11.26(d)].

Lemma 7.4. *Assume that (H.1), (H.4) and (H.5) hold. Then*

(i) $\text{prox}_{\gamma V}^{M_\gamma}$ is single-valued and Lipschitz continuous with constant $(1 - \gamma(r + 2\|\mathbf{X}\|^2/\beta))^{-1}$. Moreover, $\text{prox}_{\gamma V}^{M_\gamma} \circ M_\gamma^{-1}$ is monotone.

(ii) $M_{\gamma, \gamma} V$ is $C^1(\mathbb{R}^M)$ and its gradient is Lipschitz-continuous and hypomonotone.

Hypomonotonicity of the drift coefficient has been exploited for backward (implicit) Euler discretization of SDEs in [31] (see also references therein).

7.2.3 Algorithm and guarantees

We are now ready to state our scheme to sample from μ_n . Define the following SDE

$$d\mathbf{L}(t) = -\frac{1}{2}\nabla\left(M_{\gamma, \gamma} V \circ M_\gamma^{-1/2}\right)(\mathbf{L}(t))dt + M_\gamma^{1/2}d\mathbf{W}(t), \quad t > 0. \quad (7.4)$$

where G fulfills (H.1) and (H.4), and γ obeys (H.5). Owing to Lemma 7.4, it follows from [33, Theorem 2.4-6.1, Chapter IV] that for every initial point $\mathbf{L}(0)$ such that $\mathbb{E}\left[\|\mathbf{L}(0)\|_2^2\right] < \infty$, there exists a unique solution to the SDE (7.4) which is strongly Markovian, and the diffusion is non-explosive, i.e. $\mathbb{E}\left[\|\mathbf{L}(t)\|_2^2\right] < \infty$ for all $t > 0$. Moreover, by [49, Theorem 2.1], \mathbf{L} admits an (unique) invariant measure μ_γ having a density $\boldsymbol{\theta} \mapsto \exp\left(-M_{\gamma, \gamma} V(M_\gamma^{-1/2}\boldsymbol{\theta})\right)/Z_\gamma$, where

$$Z_\gamma = \sqrt{\det(M_\gamma)} \int_{\mathbb{R}^M} \exp\left(-M_{\gamma, \gamma} V(\mathbf{u})\right) d\mathbf{u}.$$

The following proposition answers the natural question on the behaviour of $\mu_\gamma - \mu_n$ as a function of γ .

Proposition 7.1. *Assume that (H.1), (H.4) and (H.5) hold. Then, μ_γ converges to μ_n in total variation as $\gamma \rightarrow 0$.*

In fact, the proof of this proposition only needs (H.1) and (H.2). Assumptions (H.4) and (H.5) are required only for (7.4) to be well-posed.

Inserting the identities of Lemma 7.1 and 7.3 into (7.4), and applying the change of variable $\mathbf{U}(t) = M_\gamma^{-1/2}\mathbf{L}(t)$, we get the SDE

$$d\mathbf{U}(t) = -\frac{1}{2}\left(\frac{\mathbf{I}_M - \text{prox}_{\gamma G} \circ (\mathbf{I}_M - \gamma\nabla F)}{\gamma}\right)(\mathbf{U}(t))dt + d\mathbf{W}(t), \quad \mathbf{U}(0) = \mathbf{u}_0, \quad t > 0. \quad (7.5)$$

With the same arguments as above, this SDE is well-posed (has a unique strong solution and the diffusion is non-explosive).

Consider now the forward Euler discretization of (7.5) with step-size $\delta > 0$, which can be rearranged as

$$\mathbf{U}_{k+1} = \left(1 - \frac{\delta}{2\gamma}\right)\mathbf{U}_k + \frac{\delta}{2\gamma}\text{prox}_{\gamma G}(\mathbf{U}_k - \gamma\nabla F(\mathbf{U}_k)) + \sqrt{\delta}\mathbf{Z}_k, \quad t > 0, \quad \mathbf{U}_0 = \mathbf{u}_0. \quad (7.6)$$

Clearly, without the stochastic term $\sqrt{\delta}\mathbf{Z}_k$, this would be a relaxed version of the classical Forward-Backward proximal splitting algorithm with relaxation parameter $\delta/(2\gamma)$ and descent step-size γ [4]. Hence, we coin the scheme (7.6) FB-LMC. Taking $\delta = 2\gamma$, we obtain its "unrelaxed" version. Observe again that by Lemma 7.1 and 7.3, this is also equivalent to a gradient descent on the Moreau envelope of V in the metric M_γ with step-size δ .

Remark 7.1. One can propose yet another algorithm to approximately sample from μ_n according to the following strategy. Consider smoothing G (resp. J if $G = J \circ \mathbf{D}$ as for analysis-type prior as in (5.2)) by replacing it with its Moreau envelope $\mathbf{I}_{M,\gamma}G$ (resp. with $\mathbf{I}_{M,\gamma}J \circ \mathbf{D}$), where $\gamma \in]0, 1/r[$. In this case Lemma 7.4 applies to G (resp. J), whence we deduce that $\mathbf{I}_{M,\gamma}G$ (resp. $\mathbf{I}_{M,\gamma}J \circ \mathbf{D}$) is $C^1(\mathbb{R}^M)$ with Lipschitz-continuous and hypomonotone gradient. Proposition 7.1 can also be modified correspondingly to show that $\pi_\gamma^G(d\theta) \propto \exp(-\mathbf{I}_{M,\gamma}G)(\theta)$ converges in total variation to π as $\gamma \rightarrow 0$. Thus one can replace \mathbf{M}_γ by \mathbf{I}_M in (7.4), and $\nabla(\mathbf{M}_{\gamma,\gamma}V \circ \mathbf{M}_\gamma^{-1/2})$ with $\nabla F + \nabla \mathbf{I}_{M,\gamma}G = \nabla F + \gamma^{-1}(\mathbf{I}_M - \text{prox}_{\gamma G})$ (resp. $\nabla F + \gamma^{-1}\mathbf{D}^T \circ (\mathbf{I}_M - \text{prox}_{\gamma J}) \circ \mathbf{D}$). This is a generalization beyond the convex case to the scheme described in [26]. We do not pursue this further here for space limitation and leave the details to a forthcoming work.

From (7.6), an Euler approximate solution is defined as

$$\mathbf{U}^\delta(t) \stackrel{\text{def}}{=} \mathbf{U}_0 - \frac{1}{2\gamma} \int_0^t (\bar{\mathbf{U}}(s) - \text{prox}_{\gamma G} \circ (\bar{\mathbf{U}}(s) - \gamma \nabla F(\bar{\mathbf{U}}(s)))) ds + \int_0^t d\mathbf{W}(s) ds,$$

where $\bar{\mathbf{U}}(t) = \mathbf{U}_k$ for $t \in [k\delta, (k+1)\delta[$. Observe that $\mathbf{U}^\delta(k\delta) = \bar{\mathbf{U}}(k\delta) = \mathbf{U}_k$, hence $\mathbf{U}^\delta(t)$ and $\bar{\mathbf{U}}(t)$ are continuous-time extensions to the discrete-time chain $\{\mathbf{U}_k\}_k$.

Mean square convergence of the pathwise approximation (7.6) and of its first-order moment can be established as follows.

Theorem 7.1. Assume that (H.1), (H.4) and (H.5) hold. Then

$$\|\mathbb{E}[\mathbf{U}^\delta(T)] - \mathbb{E}[\mathbf{U}(T)]\|_2 \leq \mathbb{E} \left[\sup_{0 \leq t \leq T} \|\mathbf{U}^\delta(t) - \mathbf{U}(t)\|_2 \right] = O(\delta^{1/2}) \xrightarrow{\delta \rightarrow 0} 0.$$

7.3 Application to group analysis sparsity prior

Let's now consider the prior (5.2) with g chosen according to Example 5.2. For the sake of simplicity, we take $R = +\infty$ (implying $H = +\infty$ by Proposition 5.1)². Thus

$$G(\theta) = \sum_{l=1}^L \left[\alpha^a \|[D\theta]_{\mathcal{G}_l}\|_2^a + c \log \left(\tau^b + \|[D\theta]_{\mathcal{G}_l}\|_2^b \right) \right],$$

and G is proper, continuous and bounded below, i.e. it fullfils (H.1). This function also has a supporting quadratic minorant at each point, and thus verifies (H.4).

It remains now to compute the proximal mapping of G . Owing to Assumption 4.2, one can show quite immediately that

$$\text{prox}_{\gamma G} = \widetilde{\mathbf{D}}^T \circ \text{prox}_{\gamma J + \iota_{\text{Im}(\mathbf{D})}} \circ \mathbf{D},$$

where $G = J \circ \mathbf{D}$, $\iota_{\text{Im}(\mathbf{D})}(\mathbf{u}) = 0$ if $\mathbf{u} \in \text{Im}(\mathbf{D})$ and $+\infty$ otherwise, and $\widetilde{\mathbf{D}}\widetilde{\mathbf{D}}^T$ is indeed definite positive on $\text{Im}(\mathbf{D})$. This means that, in general, G does not have an easy-to-compute expression, but rather necessitates to solve an optimization subproblem. To overcome this difficulty, we consider the case where \mathbf{D} is invertible. Observe that this assumption can be removed using the algorithmic alternative outlined in Remark 7.1. In

²The exposition can be generalized to the case of bounded R , in which case G will involve the indicator function of Θ , and still verifies (H.1).

this case, one can operate a simple change of variable $\mathbf{u} = \mathbf{D}\boldsymbol{\theta}$ and replace F with $F \circ \mathbf{D}^{-1}$, and G with $J = G \circ \mathbf{D}^{-1}$. The latter has a closed form expression as shown in the following two lemmas.

First, observe that J is separable into L independent functions (one for each group \mathcal{G}_l), and each of these functions depends only on the ℓ_2 norm of the block, whence the proximal mapping has the useful form.

Lemma 7.5. *For any $\mathbf{u} \in \mathbb{R}^P$ and $\gamma > 0$, we have*

$$\text{prox}_{\gamma J}(\mathbf{u}) = \begin{pmatrix} \text{prox}_{\gamma w}(\|\mathbf{u}_{\mathcal{G}_1}\|_2) \frac{\mathbf{u}_{\mathcal{G}_1}}{\|\mathbf{u}_{\mathcal{G}_1}\|_2} \\ \vdots \\ \text{prox}_{\gamma w}(\|\mathbf{u}_{\mathcal{G}_P}\|_2) \frac{\mathbf{u}_{\mathcal{G}_P}}{\|\mathbf{u}_{\mathcal{G}_P}\|_2} \end{pmatrix}$$

where $w : t \in \mathbb{R}^+ \mapsto w(t) = \alpha t^a + c \log(\tau^b + t^b)$.

Remark 7.2. *Lemma 7.5 remains valid for any prior that depends only on \mathbf{u} through the norms of its groups, e.g. (5.2), in which case $w(t) = \alpha t^a - \log(g(t))$, with the proviso that $-\log \circ g$ is lsc.*

We now turn to the computation of $\text{prox}_{\gamma w}$. This problem is treated in the following lemma.

Lemma 7.6. *Let w as defined in Lemma 7.5 and $\gamma > 0$. Then*

$$\text{prox}_{\gamma w}(t) = \begin{cases} 0 & \text{if } |t| \leq p_0, \\ t - \text{sgn}(t)\gamma w'(|\text{prox}_{\gamma w}(t)|) & \text{if } |t| > p_0, \end{cases} \quad (7.7)$$

where $p_0 = \min_{t \geq 0} t + \gamma w'(t)$.

Clearly, the proximal mapping of $J = G \circ \mathbf{D}^{-1}$ is single-valued hence conforming with (H.3). In addition, computing it boils down to solving the simple one-dimensional non-linear equation (7.7), which can be done very efficiently.

8 Numerical experiments

In this section, some numerical experiments are conducted to illustrate and validate the numerical performance of the proposed EWA estimator. We consider a linear regression problem

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\theta}_0 + \boldsymbol{\xi},$$

where $\boldsymbol{\xi}$ is the noise, \mathbf{X} is the design matrix. $\boldsymbol{\theta}_0 \in \mathbb{R}^M$ is the unknown regression vector of interest assumed to obey the group-analysis sparsity assumption (Assumption 4.1), with a given invertible dictionary \mathbf{D} . We take the prior (5.2) where g is given in Example 5.2 with $H = +\infty$. We then have to choose a distribution on the noise $\boldsymbol{\xi}$ such that β is independent of H . Such type of distributions is specified in [20, Section 2]. For our implementation, we assume $\boldsymbol{\xi}$ i.i.d. $\sim \mathcal{N}(0, \sigma^2)$. The noise level σ is chosen according to the simulated $\boldsymbol{\theta}_0$.

The main steps of FB-LM applied to compute EWA with our prior are summarized in Algorithm 8.1. It takes as inputs the vector of observations \mathbf{y} , the design matrix \mathbf{X} , the noise level σ , the analysis operator \mathbf{D} and the group size G . The parameters of EWA were chosen as prescribed in our theoretical analysis. For instance, the temperature parameter is set to $\beta = 4\sigma^2$, the parameters of the prior: $a \in]0, 1]$, $b \in]0, 1]$, $c > (2 + G)/b$, $\tau \sim 1/(Mn)$ and $\alpha \leq 1/(3\sqrt{K_{a,g}L})^{1/a}$, where $K_{a,g}$ is given in Corollary 6.3. The number

of iterations N and the step-size δ are chosen respectively large and small enough to guarantee convergence and discretization consistency of the algorithm.

Algorithm 8.1 Pseudo-code of the FB-LMC algorithm.

Input: The given data $(\mathbf{y}, \mathbf{X}, \sigma, \mathbf{D}, G)$ and the parameters $(a, b, c, \alpha, \beta, \tau, \delta, T)$
Output: The vector $\hat{\boldsymbol{\theta}}_n \in \mathbb{R}^M$.
Pre-computations: $P = \text{size}(\mathbf{D}, 1); N = \text{floor}(T/\delta); \text{invD} = \text{inv}(\mathbf{D}); \mathbf{A} = \mathbf{X} * \text{invD}; \mathbf{AA} = \mathbf{A}' * \mathbf{A}; \mathbf{Ay} = \mathbf{A}' * \mathbf{y};$
Initialization: $\mathbf{Li} = \text{zeros}(P, 1); \mathbf{u_bar} = \text{zeros}(P, 1);$
for all i **in** $1:1:N$ **do**
 $\text{nablaFD} = (2/\beta) * (\mathbf{AA} * \mathbf{Li} - \mathbf{Ay});$
 $\text{epsilon} = \text{sqrt}(\delta) * \text{randn}(P, 1);$
 $\mathbf{Li} = (1 - \delta/(2\gamma)) + \delta/(2\gamma) * \text{prox}_{\gamma G \circ \mathbf{D}^{-1}}(\mathbf{Li} - \gamma * \text{nablaFD} + \text{epsilon});$
 $\mathbf{u_bar} \leftarrow \mathbf{u_bar} + \mathbf{Li};$
end for
return $\text{invD} * \mathbf{u_bar} / N;$

8.1 1-D signal recovery under group sparsity

In this example, we set $\mathbf{D} = \mathbf{I}_M$, which corresponds to the classical group sparsity. The design matrix is drawn uniformly at random from the Rademacher ensemble, i.e. its entries are i.i.d. variates valued in $\{-1, 1\}$ with equal probabilities. The non-zero entries of $\boldsymbol{\theta}_0$ are equal to 1 and we denote $S = \|\boldsymbol{\theta}_0\|_0$ the sparsity level of $\boldsymbol{\theta}_0$. Two types of sparsity behavior are considered: individual sparsity where $G_{\boldsymbol{\theta}_0} = 1$; group structured sparsity with $G_{\boldsymbol{\theta}_0} = 4$. Besides, the positions of the non-zero/active entries (for $G_{\boldsymbol{\theta}_0} = 1$) or groups (for $G_{\boldsymbol{\theta}_0} = 4$) are chosen randomly uniformly on $\{1, \dots, M\}$.

The experiments are performed by fixing $M = 128$, and taking $S = 4, 8, \dots, M, n = 8, 16, \dots, M$, step-size $\delta = 4\sigma^2/(nM)$ and integration time $T = 3500$. The parameters in the prior are chosen to minimize the remainder term in the oracle inequality (6.4). For each (S, n) , and each value of $G_{\boldsymbol{\theta}_0}$, $N_{\text{rep}} = 20$ instances of the problem suite $(\mathbf{X}, \boldsymbol{\theta}_0, \mathbf{Y})$ are generated, and EWA is applied with a chosen G and the other parameters as detailed above. The estimation quality/success is then assessed by

$$\pi_{S,n} = \frac{1}{N_{\text{rep}}} \sum_{j=1}^{N_{\text{rep}}} I\left(\|\hat{\boldsymbol{\theta}}_n^{(j,S,n)} - \boldsymbol{\theta}_0^{(j,S,n)}\|_n \leq \epsilon\right), \quad (8.1)$$

where $\epsilon > 0$ (we choose $\epsilon = 0.4$) and $\hat{\boldsymbol{\theta}}_n^{(j,S,n)}$ (resp. $\boldsymbol{\theta}_0^{(j,S,n)}$) corresponds to $\hat{\boldsymbol{\theta}}_n$ (resp. $\boldsymbol{\theta}_0$) in the j -th replication of (S, n) .

S/M and n/M are respectively normalized measures of sparsity and problem indeterminacy. We get a two-dimensional phase space $(S/M, n/M) \in [0, 1]^2$ describing the difficulty of a problem instance, i.e. problems are easier as one moves up (more measurements) and to the left (sparser $\boldsymbol{\theta}_0$). Phase diagrams plotting $\pi_{S,n}$ in (8.1) as a function $(S/M, n/M)$ were widely advocated by Donoho and co-authors for ℓ_1 minimization [21]. Such diagrams often have an interesting two-phase structure (as displayed in Figures 1(a)-(d), brighter color indicate better success), with phases separated by a specific curve, called phase transition curve. Thus, a good estimator is intended to have a large bright area which indicates its good performance at a wider range of (S, n) .

Figure 1(a) (resp. (b)) shows the phase diagrams when $G_{\theta_0} = 1$ and $G = 1$ (resp. $G = 4$) in EWA. In this case, the phase transition curve for $G = 1$, the correct group size, is slightly better than with $G = 4$. The situation reverses for Figures 1(c)-(d) where $G_{\theta_0} = 4$, and one observes that the success area is significantly better using $G = 4$ than $G = 1$. This is expected as it reveals better performance of EWA when used with the choice $G = G_{\theta_0}$. This is also confirmed by visual inspection of Figures 1(c')-(d'), where we plotted instances of recovered vectors $\hat{\theta}_n^{(j,S,n)}$ when $(S, G_{\theta_0}) \in \{4, 8\} \times \{1, 4\}$ and $n/M = 1/2$. EWA was again applied with $G = 1$ and $G = 4$ in each case. Large spurious entries appear outside the true support when the group size is not correctly chosen, though the impact is less important for $G = 1$.

It is worth observing that $S/M = \|\theta_0\|_{0,\mathcal{G}} G_{\theta_0}/M$. As far as the expected phase transition curve is concerned, one has from the oracle inequality (6.4) and Remark 6.2 that it is expected to occur for

$$n/M = C_\epsilon \|\theta_0\|_{0,\mathcal{G}} G_{\theta_0}/M (\log(M/G_{\theta_0})/G_{\theta_0}) = C_\epsilon S/M (\log(M/G_{\theta_0})/G_{\theta_0})$$

for some constant $C_\epsilon > 0$ depending on ϵ . That is, the phase transition curve is linear (M and G_{θ_0} are fixed for each diagram), which is confirmed by visual inspection of Figures 1(a)-(d), where the overlaid blue line is the fitted linear phase transition curve.

8.2 2-D image recovery under analysis group-sparsity

In the second numerical experiment, θ_0 is a 2-D image which is a matrix in $\mathbb{R}^{160 \times 160}$ (a close-up of the known Shepp-Logan phantom, see Figure 2(a)). Let us denote vec the vectorization operator. Thus $\text{vec}(\theta_0)$ is vector in \mathbb{R}^M with $M = 160^2$, and our goal is to recover θ_0 from

$$Y = X \text{vec}(\theta_0) + \xi,$$

where $\xi \sim \mathcal{N}(0, \sigma^2 I_n)$ and $X \in \mathbb{R}^{n \times M}$ is again random whose entries are i.i.d. from the Rademacher distribution. Since the targeted image is piecewise-constant, a popular prior is the so-called isotropic total variation [51]. It turns out that this can be cast in our analysis group-sparsity framework as a special case.

Let $D_1 : \mathbb{R}^{\sqrt{M} \times \sqrt{M}} \rightarrow \mathbb{R}^{\sqrt{M} \times \sqrt{M}}$ and $D_2 : \mathbb{R}^{\sqrt{M} \times \sqrt{M}} \rightarrow \mathbb{R}^{\sqrt{M} \times \sqrt{M}}$ be the finite difference operators along, respectively, the columns and rows, with appropriate boundary conditions that we will specify shortly. We define the linear analysis operator D as

$$D : \theta \in \mathbb{R}^{\sqrt{M} \times \sqrt{M}} \mapsto (D_1(\theta), D_2(\theta)) \in \mathbb{R}^{\sqrt{M} \times \sqrt{M}} \times \mathbb{R}^{\sqrt{M} \times \sqrt{M}},$$

The isotropic total variation prior on θ promotes sparsity of $\text{vec} \left(\left\| [D(\theta)]_{i,j} \right\|_2 \right)_{1 \leq i,j \leq \sqrt{M}}$. By defining the set of groups by

$$\mathcal{G} = \bigcup_{(i,j) \in \{1, \dots, s_{p_0}\}^2} \{(i, j, 1), (i, j, 2)\},$$

one immediately realizes that measuring sparsity of the above vectorized form is equivalent to group sparsity of $D(\theta_0)$ with groups of size 2 along the third dimension. One can then use FB-LMC Algorithm 8.1 to get the EWA estimator. We have however to ensure that D_1 and D_2 are invertible, which is achieved by a proper choice of the boundary conditions (not periodic of course).

The results are depicted in Figure 2. The number of observations is $n = 9/16M = 14400$, and we have $\|D(\theta_0)\|_{0,\mathcal{G}} = 1376 \ll n$. A notable property of the EWA estimate is that it does not suffer from the stair-casing effect, unlike total variation minimization.

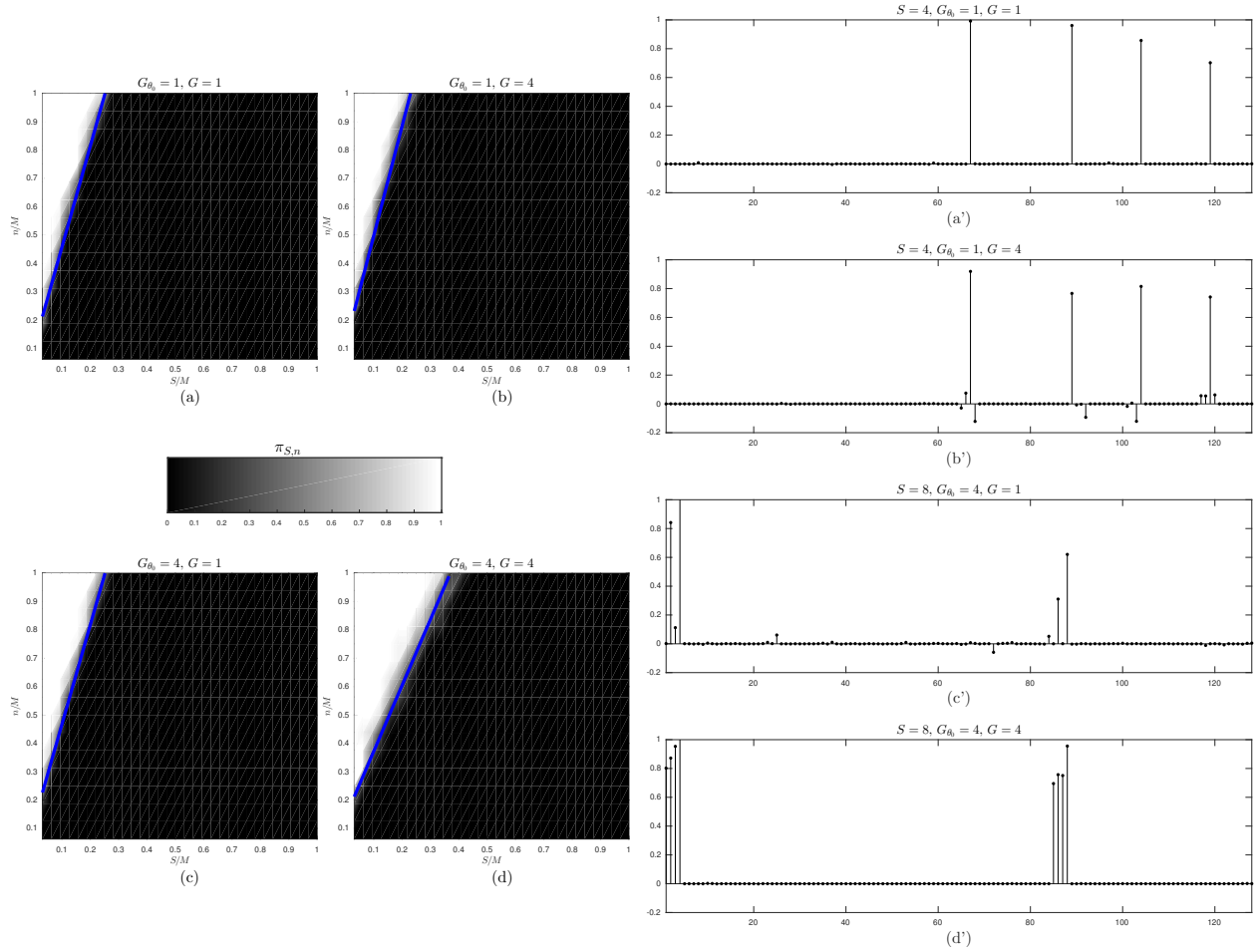


Figure 1: (a)-(d): Phase diagrams of EWA for $D = I_M$, the color bar ranges from dark ($\pi_{S,n} = 0$) to bright ($\pi_{S,n} = 1$). The blue line is the fitted phase transition curve. (a')-(d'): Examples of vectors $\hat{\theta}_n^{(j,S,n)}$ recovered by EWA with $n/M = 1/2$, two sparsity levels $S = 4$ and $S = 8$ and two group sizes $G_{\theta_0} = 1$ and $G_{\theta_0} = 4$.

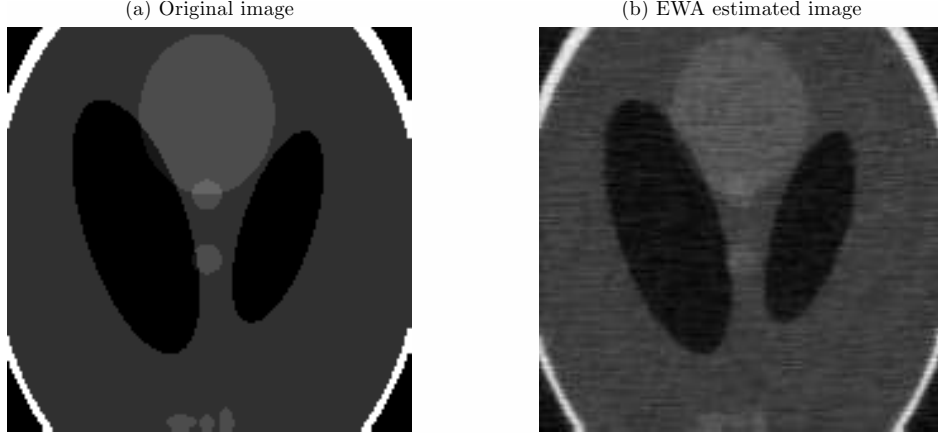


Figure 2: (a): Original close-up of Shepp-Logan phantom image. (b): Image recovered by EWA with $\delta = 2 \cdot 10^{-8}$ and $T = 10^4$.

9 Conclusion

In this paper, we proposed a class of EWA estimators constructed from a novel and versatile family of priors which promotes analysis group-sparsity, where the analysis operator corresponds to a frame. Its quality is guaranteed by establishing a sharp SOI with a small remainder term in high-dimension. We also described a forward-backward proximal LMC algorithm, which is an implementation of EWA and can be viewed as a Forward Euler discretization of a Langevin diffusion involving the Moreau-Envelope of the potential in a proper metric. We derived consistency guarantees of this discretization. The performance of the estimator was illustrated on some numerical experiments which support our theoretical findings. There are still open problems that we leave to a future work. More precisely, one direction is to investigate how to remove the frame assumption. Another one would be to derive further/better quantitative convergence bounds of the proposed discretization to the target distribution and/or to relax some of the assumptions to cover an even larger class of distributions.

10 Proofs

10.1 Proofs of auxiliary results

10.1.1 Proof of Lemma 4.1

Consider the linear mapping $\mathbf{D} : \boldsymbol{\theta} \in \mathbb{R}^M \mapsto \mathbf{D}\boldsymbol{\theta} \in \mathbb{R}^P$, $P \geq M$. The Jacobian matrix of this mapping is obviously \mathbf{D} for any $\boldsymbol{\theta} \in \mathbb{R}^M$. Since \mathbf{D} is a frame, it is injective, hence so-called M -regular (see [41, Section 1.5]). In particular, $\det(\mathbf{D}^T \mathbf{D}) > 0$. Thus combining [41, Theorems 1.12 and 3.4] and the Cauchy-Binet formula [41, Theorem 3.3]), we have the change of variables formula

$$\begin{aligned} \int_{\Theta} u(\mathbf{D}\boldsymbol{\theta}) d\boldsymbol{\theta} &= \frac{1}{\sqrt{\det(\mathbf{D}^T \mathbf{D})}} \int_{\mathbb{R}^P} \left(\sum_{\boldsymbol{\theta} \in \Theta \cap \{\boldsymbol{\omega} : \mathbf{D}\boldsymbol{\omega} = \mathbf{v}\}} u(\mathbf{D}\boldsymbol{\theta}) \right) d\mathbf{v} \\ &= \frac{1}{\sqrt{\det(\mathbf{D}^T \mathbf{D})}} \int_{\text{Im}(\mathbf{D})} \left(\sum_{\boldsymbol{\theta} \in \Theta \cap \{\boldsymbol{\omega} : \mathbf{D}\boldsymbol{\omega} = \mathbf{v}\}} u(\mathbf{D}\boldsymbol{\theta}) \right) d\mathbf{v}. \end{aligned} \quad (10.1)$$

Using once again that \mathbf{D} is a frame, i.e. it is bijective on its image $\text{Im}(\mathbf{D})$, the result follows. This concludes the proof. \square

10.1.2 Proof of Proposition 5.1

Let $\boldsymbol{\theta} \in \Theta$ and $i \in \{1, \dots, n\}$. Setting $u_i^\theta = \sum_{j=1}^M \theta_j f_j(x_i)$ and $\mathbf{u}^\theta = (u_1^\theta, \dots, u_n^\theta)^T$, and by virtue of (4.1), we have

$$\|\mathbf{u}^\theta\|_2 = \|\mathbf{X}\boldsymbol{\theta}\|_2 \leq \|\mathbf{X}\| \|\widetilde{\mathbf{D}}^T\| \|\mathbf{D}\boldsymbol{\theta}\|_2. \quad (10.2)$$

From (4.2), we get also

$$\|\|\widetilde{\mathbf{D}}^T\|\| = \sqrt{\sigma_1(\widetilde{\mathbf{D}}^T \widetilde{\mathbf{D}}^T)} \leq \frac{1}{\sqrt{\mu}}. \quad (10.3)$$

Moreover, since $a \in]0, 1]$, by the triangle inequality, we get

$$\|\mathbf{D}\boldsymbol{\theta}\|_2 \leq \sum_{l=1}^L \|\mathbf{D}\boldsymbol{\theta}\|_{\mathcal{G}_l} \leq \left[\sum_{l=1}^L \|\mathbf{D}\boldsymbol{\theta}\|_{\mathcal{G}_l}^a \right]^{1/a} = \|\mathbf{D}\boldsymbol{\theta}\|_{a, \mathcal{G}} \leq R^{1/a}. \quad (10.4)$$

Combining (10.2), (10.3) and (10.4), we obtain

$$\|\mathbf{u}^\theta\|_2 \leq \frac{\|\mathbf{X}\| R^{1/a}}{\sqrt{\mu}},$$

which in turn implies $\mathbf{u}^\theta \in \mathcal{B}$. Therefore, for any $(\boldsymbol{\theta}, \boldsymbol{\theta}') \in \Theta^2$,

$$\|\mathbf{f}_\theta - \mathbf{f}_{\theta'}\|_2 = \|\mathbf{L}(\mathbf{u}^\theta) - \mathbf{L}(\mathbf{u}^{\theta'})\|_2 \leq 2 \max_{\mathbf{x} \in \mathcal{B}} \|\mathbf{L}(\mathbf{x})\|_2.$$

\square

10.1.3 Proof of Lemma 5.1

Let us first check the integrability condition (5.3). By Lemmas 4.1 and 2.2, we obtain

$$\begin{aligned} \int_{\mathbb{R}^M} \prod_{l=1}^L g(\|\mathbf{D}\mathbf{u}\|_{\mathcal{G}_l}) d\mathbf{u} &= \frac{1}{\sqrt{\det(\mathbf{D}^T \mathbf{D})}} \int_{\text{Im}(\mathbf{D})} \prod_{l=1}^L g(\|\mathbf{v}_{\mathcal{G}_l}\|_2) d\mathbf{v} \\ &\leq \frac{1}{\sqrt{\det(\mathbf{D}^T \mathbf{D})}} \int_{\mathbb{R}^P} \prod_{l=1}^L g(\|\mathbf{v}_{\mathcal{G}_l}\|_2) d\mathbf{v} = \frac{1}{\sqrt{\det(\mathbf{D}^T \mathbf{D})}} \left(\int_{\mathbb{R}^G} g(\|\mathbf{u}\|_2) d\mathbf{u} \right)^L \\ &= \frac{C_G^L}{\sqrt{\det(\mathbf{D}^T \mathbf{D})}} \left(\int_0^\infty z^{G-1} g(z) dz \right)^L. \end{aligned}$$

Since $G \geq 1$ and $g : \mathbb{R}^+ \rightarrow \mathbb{R}^+$, by (5.6), we get

$$\begin{aligned}
& \int_{\mathbb{R}^M} \prod_{l=1}^L g(\|[\mathbf{D}\mathbf{u}]_{\mathcal{G}_l}\|_2) d\mathbf{u} \\
& \leq \frac{C_G^L}{\sqrt{\det(\mathbf{D}^T \mathbf{D})}} \left(\int_0^1 z^{G-1} g(z) dz + \int_1^\infty w^{G-1} g(w) dw \right)^L \\
& \leq \frac{C_G^L}{\sqrt{\det(\mathbf{D}^T \mathbf{D})}} \left(\sup_{z \in [0,1]} g(z) + \int_1^\infty w^{G+1} g(w) dw \right)^L < \infty.
\end{aligned} \tag{10.5}$$

Therefore, g satisfies Assumption 5.2-2.

Now, we check the moment condition (5.4). Using similar arguments to the bound (10.5), we have

$$\begin{aligned}
& \int_{\mathbb{R}^M} \|\mathbf{D}\mathbf{u}\|_2^2 \prod_{k=1}^L g(\|[\mathbf{D}\mathbf{u}]_{\mathcal{G}_k}\|_2) d\mathbf{u} \\
& \leq \frac{1}{\sqrt{\det(\mathbf{D}^T \mathbf{D})}} \int_{\mathbb{R}^P} \|\mathbf{v}_{\mathcal{G}_i}\|_2^2 \prod_{k=1}^L g(\|\mathbf{v}_{\mathcal{G}_k}\|_2) d\mathbf{v} \\
& = \frac{1}{\sqrt{\det(\mathbf{D}^T \mathbf{D})}} \left(\int_{\mathbb{R}^G} \|\mathbf{u}\|_2^2 g(\|\mathbf{u}\|_2) d\mathbf{u} \right) \left(\int_{\mathbb{R}^G} g(\|\mathbf{v}\|_2) d\mathbf{v} \right)^{L-1} \\
& = \frac{C_G^L}{\sqrt{\det(\mathbf{D}^T \mathbf{D})}} \int_0^\infty z^{G+1} g(z) dz \left(\int_0^\infty w^{G-1} g(w) dw \right)^{L-1} < \infty,
\end{aligned} \tag{10.6}$$

whence we conclude that g satisfies Assumption 5.2-3. \square

10.1.4 Proof of Corollary 6.1

Let $\gamma = 3$, $\nu = 2$ and $\eta = 1$. We have $\gamma/\nu < \eta + 1$ so that Lemma 2.1 applies. We thus obtain

$$\int_0^\infty z^{G+1} g(z) dz = \int_0^\infty \frac{z^2}{(z^2 + \tau^2)^2} dz < \infty.$$

From Lemma 5.1, g satisfies Assumptions 5.2-2 and 5.2-3. Moreover, taking $h(t) = 1 + t/\tau$ and $\lambda = 4$, for all $(t, t^*) \in \mathbb{R}^2$, we have by Young's inequality

$$\begin{aligned}
\frac{g(|t - t^*|)}{g(|t|)} & = \left[\frac{\tau^2 + t^2}{\tau^2 + (t - t^*)^2} \right]^2 = \left[1 + \frac{2\tau(t - t^*)t^*/\tau + t^{*2}}{\tau^2 + (t - t^*)^2} \right]^2 \\
& \leq \left[1 + \frac{|t^*|}{\tau} + \frac{t^{*2}}{\tau^2} \right]^2 \leq \left[1 + \frac{|t^*|}{\tau} \right]^4 = h(|t^*|)^\lambda.
\end{aligned}$$

Therefore, g satisfies Assumption 5.2 for $G = 1$.

Owing to Remark 5.2 and Lemma 2.1, we obtain

$$K_{1,g} = \frac{\int_0^\infty \frac{x^2}{(\tau^2+x^2)^2} dx}{\int_0^\infty \frac{1}{(\tau^2+y^2)^2} dy} = \tau^2.$$

We are now in position to apply Theorem 6.1 with $\mathbf{D} = \mathbf{I}_M$ (then $P = M$), $G = 1$ (then $L = P$), $a = 1$ and $\alpha \leq 1/(3M\tau)$ to conclude. This completes the proof. \square

10.1.5 Proof of Corollary 6.2

Let $\gamma = 2 + G$, $\nu = b$ and $\eta = c - 1$. We have $\gamma/\nu < \eta + 1$ and thus Lemma 2.1 applies, whence we obtain

$$\int_0^\infty x^{G+1} g(x) dx = \int_0^\infty \frac{x^{G+1}}{(\tau^b + x^b)^c} dx < \infty.$$

From Lemma 5.1, g satisfies Assumptions 5.2-2 and 5.2-3. Recall that $b \in]0, 1]$. Taking $h(x) = 1 + (x/\tau)^b$ and $\lambda = c$, for all $(\mathbf{t}, \mathbf{t}^*) \in \mathbb{R}^G \times \mathbb{R}^G$, we have

$$\begin{aligned} \frac{g(\|\mathbf{t} - \mathbf{t}^*\|_2)}{g(\|\mathbf{t}\|_2)} &= \left[\frac{\tau^b + \|\mathbf{t}\|_2^b}{\tau^b + \|\mathbf{t} - \mathbf{t}^*\|_2^b} \right]^c \leq \left[\frac{\tau^b + \|\mathbf{t} - \mathbf{t}^*\|_2^b + \|\mathbf{t}^*\|_2^b}{\tau^b + \|\mathbf{t} - \mathbf{t}^*\|_2^b} \right]^c \\ &\leq \left[1 + \frac{\|\mathbf{t}^*\|_2^b}{\tau^b + \|\mathbf{t} - \mathbf{t}^*\|_2^b} \right]^c \leq \left[1 + \left(\|\mathbf{t}^*\|_2 / \tau \right)^b \right]^c. \end{aligned}$$

Therefore, g satisfies Assumption 5.2 with any $G \geq 1$. Applying Theorem 6.1, we obtain the claimed SOI. The proof is complete. \square

10.1.6 Proof of Corollary 6.3

Since g satisfies Assumptions 5.2-2, 5.2-3 and \mathbf{D} is invertible, by Remark 5.2 and Lemma 2.1, we get

$$K_{a,g} = \frac{\int_0^\infty \frac{r^{G-1+2a}}{(\tau^b+r^b)^c} dr}{\int_0^\infty \frac{q^{G-1}}{(\tau^b+q^b)^c} dq} = \tau^{2a} \frac{\Gamma(\frac{2a+G}{b})\Gamma(c - \frac{2a+G}{b})}{\Gamma(\frac{G}{b})\Gamma(c - \frac{G}{b})} = \tilde{K}_{a,g} \tau^{2a}.$$

According to (6.3) and the fact that $\alpha \leq 1/\left(3\tau^a \sqrt{\tilde{K}_{a,g}L}\right)^{1/a}$
 $= 1/(3\sqrt{K_{a,g}L})^{1/a}$, we obtain (6.4). This ends the proof. \square

10.2 Proofs of the SOI results

10.2.1 Proof of Theorem 6.1

Remind the prior $\pi(d\boldsymbol{\theta})$ from (5.2), where $\Theta = \left\{ \boldsymbol{\theta} \in \mathbb{R}^M : \|\mathbf{D}\boldsymbol{\theta}\|_{a,\mathcal{G}}^a \leq R \right\}$.

Let $r_L = 3\sqrt{K_{a,g}L}$, $\Theta_{p_0^D} = \left\{ \boldsymbol{\theta} \in \mathbb{R}^M : \|\mathbf{D}\boldsymbol{\theta} - \mathbf{D}\boldsymbol{\theta}^*\|_{a,\mathcal{G}}^a \leq r_L \right\}$ and

$$\boldsymbol{\theta}^* \in \left\{ \boldsymbol{\theta} \in \mathbb{R}^M : \|\mathbf{D}\boldsymbol{\theta}\|_{a,\mathcal{G}}^a \leq R - 3\sqrt{K_{a,g}L} = R - r_L \right\}. \quad (10.7)$$

We define the probability measure

$$p_0^D(d\boldsymbol{\theta}) = \frac{1}{C_L} \left(\frac{d\pi}{d\boldsymbol{\theta}}(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \right) I_{\Theta_{p_0^D}}(\boldsymbol{\theta}) d\boldsymbol{\theta},$$

where $C_L > 0$ is the normalization factor for p_0^D .

Since $r_L < R$, $\boldsymbol{\theta} \in \Theta_{p_0^D}$ implies that $\boldsymbol{\theta} - \boldsymbol{\theta}^* \in \Theta$. Therefore,

$$\begin{aligned} p_0^D(d\boldsymbol{\theta}) &= \frac{1}{C_L} \prod_{l=1}^L \exp\left(-\alpha^a \|[D\boldsymbol{\theta} - D\boldsymbol{\theta}^*]_{\mathcal{G}_l}\|_2^a\right) g\left(\|[D\boldsymbol{\theta} - D\boldsymbol{\theta}^*]_{\mathcal{G}_l}\|_2\right) \\ &\quad I_{\Theta}(\boldsymbol{\theta} - \boldsymbol{\theta}^*) I_{\Theta_{p_0^D}}(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &= \frac{1}{C_L} \prod_{l=1}^L \exp\left(-\alpha^a \|[D\boldsymbol{\theta} - D\boldsymbol{\theta}^*]_{\mathcal{G}_l}\|_2^a\right) g\left(\|[D\boldsymbol{\theta} - D\boldsymbol{\theta}^*]_{\mathcal{G}_l}\|_2\right) \\ &\quad I_{\Theta_{p_0^D}}(\boldsymbol{\theta}) d\boldsymbol{\theta}. \end{aligned}$$

For any $i \in \{1, \dots, n\}$, with $\mathbf{X}_i = (f_1(x_i), \dots, f_M(x_i))^T$, one can write

$$f_{\boldsymbol{\theta}}(x_i) = \ell\left(\sum_{j=1}^M \theta_j f_j(x_i)\right) = \ell(\mathbf{X}_i^T \boldsymbol{\theta}).$$

Taylor-Lagrange formula then gives us

$$\begin{aligned} (f_{\boldsymbol{\theta}}(x_i) - f(x_i))^2 &\leq (f_{\boldsymbol{\theta}^*}(x_i) - f(x_i))^2 + C_{f,\ell} [\mathbf{X}_i^T (\boldsymbol{\theta} - \boldsymbol{\theta}^*)]^2 \\ &\quad + 2(f_{\boldsymbol{\theta}^*}(x_i) - f(x_i)) \ell'(\mathbf{X}_i^T \boldsymbol{\theta}^*) \mathbf{X}_i^T (\boldsymbol{\theta} - \boldsymbol{\theta}^*), \end{aligned} \quad (10.8)$$

where $C_{f,\ell} = \|\ell'\|_{\infty}^2 + \|\ell''\|_{\infty} (\|\ell'\|_{\infty} + \|f\|_{\infty})$.

By summing, normalizing by $1/n$, taking the integral in Θ w.r.t. p_0^D , inequality (10.8) becomes

$$\begin{aligned} \int_{\Theta} \|f_{\boldsymbol{\theta}} - f\|_n^2 p_0^D(d\boldsymbol{\theta}) &\leq \|f_{\boldsymbol{\theta}^*} - f\|_n^2 \\ &\quad + \frac{2}{n} \sum_{i=1}^n (f_{\boldsymbol{\theta}^*}(x_i) - f(x_i)) \ell'(\mathbf{X}_i^T \boldsymbol{\theta}^*) \\ &\quad \mathbf{X}_i^T \int_{\Theta} (\boldsymbol{\theta} - \boldsymbol{\theta}^*) p_0^D(d\boldsymbol{\theta}) \\ &\quad + C_{f,\ell} \int_{\mathbb{R}^M} \frac{1}{n} \sum_{i=1}^n [\mathbf{X}_i^T (\boldsymbol{\theta} - \boldsymbol{\theta}^*)]^2 p_0^D(d\boldsymbol{\theta}). \end{aligned} \quad (10.9)$$

Note that, the right term of inequality (10.9) corresponds to a sum of three components. In the following, we keep the first component and treat the other two.

Let us first show that the second component vanishes. Indeed, let $\boldsymbol{\theta} \in \Theta_{p_0^D}$, from (10.7) and the fact that $a \in]0, 1]$, we have

$$\begin{aligned} \|D\boldsymbol{\theta}\|_{a,\mathcal{G}}^a &= \sum_{l=1}^L \|[D\boldsymbol{\theta}]_{\mathcal{G}_l}\|_2^a \leq \sum_{l=1}^L (\|[D\boldsymbol{\theta} - D\boldsymbol{\theta}^*]_{\mathcal{G}_l}\|_2 + \|[D\boldsymbol{\theta}^*]_{\mathcal{G}_l}\|_2)^a \\ &\leq \|D\boldsymbol{\theta} - D\boldsymbol{\theta}^*\|_{a,\mathcal{G}}^a + \|D\boldsymbol{\theta}^*\|_{a,\mathcal{G}}^a \leq r_L + \|D\boldsymbol{\theta}^*\|_{a,\mathcal{G}}^a \leq R. \end{aligned}$$

Then $\boldsymbol{\theta} \in \left\{ \boldsymbol{\theta} \in \mathbb{R}^M : \|\mathbf{D}\boldsymbol{\theta}\|_{a,\mathcal{G}}^a \leq R \right\} = \Theta$. Therefore, we have the embedding

$$\Theta_{p_0^{\mathcal{D}}} \subseteq \Theta. \quad (10.10)$$

In what follows, we denote $\mathbb{B}_{a,\mathcal{G}}^a(x) = \left\{ \mathbf{z} \in \mathbb{R}^P : \|\mathbf{z}\|_{a,\mathcal{G}}^a \leq x \right\}$, $\forall x > 0$ for brevity. By (10.10), property (4.1), Lemma 4.1 and symmetry of $\mathbb{B}_{a,\mathcal{G}}^a(r_L) \cap \text{Im}(\mathbf{D})$, we obtain

$$\begin{aligned} & \int_{\Theta} (\boldsymbol{\theta} - \boldsymbol{\theta}^*) p_0^{\mathcal{D}}(d\boldsymbol{\theta}) \\ & \propto \int_{\Theta \cap \Theta_{p_0^{\mathcal{D}}}} (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \prod_{l=1}^L \exp\left(-\alpha^a \|\mathbf{D}\boldsymbol{\theta} - \mathbf{D}\boldsymbol{\theta}^*\|_{\mathcal{G}_l}^a\right) g\left(\|\mathbf{D}\boldsymbol{\theta} - \mathbf{D}\boldsymbol{\theta}^*\|_{\mathcal{G}_l}\right) d\boldsymbol{\theta} \\ & = \int_{\Theta_{p_0^{\mathcal{D}}}} (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \prod_{l=1}^L \exp\left(-\alpha^a \|\mathbf{D}\boldsymbol{\theta} - \mathbf{D}\boldsymbol{\theta}^*\|_{\mathcal{G}_l}^a\right) g\left(\|\mathbf{D}\boldsymbol{\theta} - \mathbf{D}\boldsymbol{\theta}^*\|_{\mathcal{G}_l}\right) d\boldsymbol{\theta} \\ & = \frac{\widetilde{\mathbf{D}}^T}{\sqrt{\det(\mathbf{D}^T \mathbf{D})}} \int_{\mathbb{B}_{a,\mathcal{G}}^a(r_L) \cap \text{Im}(\mathbf{D})} \mathbf{z} \prod_{l=1}^L \exp\left(-\alpha^a \|\mathbf{z}_{\mathcal{G}_l}\|_2^a\right) g\left(\|\mathbf{z}_{\mathcal{G}_l}\|_2\right) d\mathbf{z} = 0, \end{aligned} \quad (10.11)$$

which is the desired claim.

We now bound the last term in the right hand side of (10.9). Define

$$p_0(d\mathbf{u}) = \frac{1}{C_L \sqrt{\det(\mathbf{D}^T \mathbf{D})}} \prod_{l=1}^L \exp\left(-\alpha^a \|\mathbf{u}_{\mathcal{G}_l}\|_2^a\right) g\left(\|\mathbf{u}_{\mathcal{G}_l}\|_2\right) I_{\text{Im}(\mathbf{D}) \cap \mathbb{B}_{a,\mathcal{G}}^a(r_L)}(\mathbf{u}) d\mathbf{u}. \quad (10.12)$$

One can see that p_0 coincides with the probability measure $p_0^{\mathcal{D}}$ on \mathbb{R}^M via a change of variables of type (4.3). So, p_0 is a probability measure on \mathbb{R}^P .

For any $i, j \in \{1, \dots, L\}$, $i \neq j$, by a change of variables, we get

$$\int_{\mathbb{R}^P} \mathbf{u}_{\mathcal{G}_i} \mathbf{u}_{\mathcal{G}_j}^T p_0(d\mathbf{u}) = - \int_{\mathbb{R}^P} \mathbf{u}_{\mathcal{G}_i} \mathbf{u}_{\mathcal{G}_j}^T p_0(d\mathbf{u}),$$

so

$$\int_{\mathbb{R}^P} \mathbf{u}_{\mathcal{G}_i} \mathbf{u}_{\mathcal{G}_j}^T p_0(d\mathbf{u}) = 0. \quad (10.13)$$

For any $j \in \{1, \dots, L\}$, as all groups have the same size, we have

$$\int_{\mathbb{R}^P} \mathbf{u}_{\mathcal{G}_j} \mathbf{u}_{\mathcal{G}_j}^T p_0(d\mathbf{u}) = \int_{\mathbb{R}^P} \mathbf{u}_{\mathcal{G}_1} \mathbf{u}_{\mathcal{G}_1}^T p_0(d\mathbf{u}). \quad (10.14)$$

For a matrix $\mathbf{A} \in \mathbb{R}^{M \times M}$ and $l \in \{1, \dots, L\}$, $[\mathbf{A}]_{\mathcal{G}_l, \mathcal{G}_l}$ stands for the restriction of \mathbf{A} to rows and columns

indexed by \mathcal{G}_l . Combining (4.1), Lemma 4.1, (10.13), (10.14) and Von Neumann's trace inequality, we obtain

$$\begin{aligned}
& \int_{\mathbb{R}^M} \frac{1}{n} \sum_{i=1}^n [\mathbf{X}_i^T(\boldsymbol{\theta} - \boldsymbol{\theta}^*)]^2 p_0^{\mathbf{D}}(d\boldsymbol{\theta}) = \frac{1}{n} \int_{\mathbb{R}^P} [\mathbf{X}\widetilde{\mathbf{D}}^T \mathbf{u}]^T \mathbf{X}\widetilde{\mathbf{D}}^T \mathbf{u} p_0(d\mathbf{u}) \\
&= \frac{1}{n} \int_{\mathbb{R}^P} \text{tr} \left(\mathbf{u}^T \widetilde{\mathbf{D}}^T \mathbf{X}^T \mathbf{X} \widetilde{\mathbf{D}}^T \mathbf{u} \right) p_0(d\mathbf{u}) \\
&= \frac{1}{n} \text{tr} \left(\left(\mathbf{X}\widetilde{\mathbf{D}}^T \right)^T \mathbf{X}\widetilde{\mathbf{D}}^T \int_{\mathbb{R}^P} \mathbf{u}\mathbf{u}^T p_0(d\mathbf{u}) \right) \\
&= \frac{1}{n} \sum_{l=1}^L \text{tr} \left(\left[\left(\mathbf{X}\widetilde{\mathbf{D}}^T \right)^T \mathbf{X}\widetilde{\mathbf{D}}^T \right]_{\mathcal{G}_l, \mathcal{G}_l} \int_{\mathbb{R}^P} \mathbf{u}_{\mathcal{G}_l} \mathbf{u}_{\mathcal{G}_l}^T p_0(d\mathbf{u}) \right) \\
&\leq \frac{1}{n} \sum_{l=1}^L \sum_{j=1}^G \sigma_j \left(\left[\left(\mathbf{X}\widetilde{\mathbf{D}}^T \right)^T \mathbf{X}\widetilde{\mathbf{D}}^T \right]_{\mathcal{G}_l, \mathcal{G}_l} \right) \sigma_j \left(\int_{\mathbb{R}^P} \mathbf{u}_{\mathcal{G}_l} \mathbf{u}_{\mathcal{G}_l}^T p_0(d\mathbf{u}) \right) \\
&\leq \frac{1}{n} \int_{\mathbb{R}^P} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 p_0(d\mathbf{u}) \sum_{l=1}^L \text{tr} \left(\left[\left(\mathbf{X}\widetilde{\mathbf{D}}^T \right)^T \mathbf{X}\widetilde{\mathbf{D}}^T \right]_{\mathcal{G}_l, \mathcal{G}_l} \right) \\
&= \frac{1}{n} \text{tr} \left(\left(\mathbf{X}\widetilde{\mathbf{D}}^T \right)^T \mathbf{X}\widetilde{\mathbf{D}}^T \right) \int_{\mathbb{R}^P} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 p_0(d\mathbf{u}). \tag{10.15}
\end{aligned}$$

Moreover, by inequality (4.2), Assumption 5.1 and Von Neumann's trace inequality, we obtain

$$\begin{aligned}
& \frac{1}{n} \text{tr} \left(\left(\mathbf{X}\widetilde{\mathbf{D}}^T \right)^T \mathbf{X}\widetilde{\mathbf{D}}^T \right) = \frac{1}{n} \text{tr} \left(\mathbf{X}^T \mathbf{X} \widetilde{\mathbf{D}}^T \widetilde{\mathbf{D}}^T \right) \\
&\leq \sum_{j=1}^M \sigma_j \left(\frac{\mathbf{X}^T \mathbf{X}}{n} \right) \sigma_j \left(\widetilde{\mathbf{D}}^T \widetilde{\mathbf{D}}^T \right) \leq \sigma_1 \left(\widetilde{\mathbf{D}}^T \widetilde{\mathbf{D}}^T \right) \sum_{j=1}^M \sigma_j \left(\frac{\mathbf{X}^T \mathbf{X}}{n} \right) \leq \frac{M}{\mu}. \tag{10.16}
\end{aligned}$$

Putting together (10.15) and (10.16), we get the bound

$$C_{f,\ell} \int_{\mathbb{R}^M} \frac{1}{n} \sum_{i=1}^n [\mathbf{X}_i^T(\boldsymbol{\theta} - \boldsymbol{\theta}^*)]^2 p_0^{\mathbf{D}}(d\boldsymbol{\theta}) \leq C_{f,\ell} \frac{M}{\mu} \int_{\mathbb{R}^P} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 p_0(d\mathbf{u}). \tag{10.17}$$

Thanks to (10.11) and (10.17), inequality (10.9) becomes

$$\int_{\Theta} \|f_{\boldsymbol{\theta}} - f\|_n^2 p_0^{\mathbf{D}}(d\boldsymbol{\theta}) \leq \|f_{\boldsymbol{\theta}^*} - f\|_n^2 + C_{f,\ell} \frac{M}{\mu} \int_{\mathbb{R}^P} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 p_0(d\mathbf{u}). \tag{10.18}$$

Now, inserting (10.18) into Theorem 3.1 (with $p = p_0^{\mathbf{D}}$), we arrive at

$$\mathbb{E} \left[\|\widehat{f}_n - f\|_n^2 \right] \leq \|f_{\boldsymbol{\theta}^*} - f\|_n^2 + C_{f,\ell} \frac{M}{\mu} \int_{\mathbb{R}^P} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 p_0(d\mathbf{u}) + \frac{\beta \text{KL}(p_0^{\mathbf{D}}, \pi)}{n}. \tag{10.19}$$

To complete the proof, it remains to bound the last two terms in the right hand side of (10.19). This is the goal of the following lemma whose proof will be detailed in Section 10.2.2.

Lemma 10.1. *Consider the same framework as the one in Theorem 6.1, we have*

$$\int_{\mathbb{R}^P} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 p_0(d\mathbf{u}) \leq 2K_{1,g} e^{rL\alpha^a}, \tag{10.20}$$

and

$$\text{KL}(p_0^{\mathbf{D}}, \pi) \leq 1 + r_L \alpha^a + \lambda \sum_{l=1}^L \log \left\{ h \left(\left\| [\mathbf{D}\boldsymbol{\theta}^*]_{\mathcal{G}_l} \right\|_2 \right) \right\} + \alpha^a \left\| \mathbf{D}\boldsymbol{\theta}^* \right\|_{a, \mathcal{G}}^a. \quad (10.21)$$

With $r_L = 3\sqrt{K_{a,g}}L$, it follows from (10.19) and Lemma 10.1 that

$$\begin{aligned} \mathbb{E} \left[\left\| \widehat{f}_n - f \right\|_n^2 \right] &\leq \left\| f_{\boldsymbol{\theta}^*} - f \right\|_n^2 + \frac{\beta}{n} \left(1 + 3\sqrt{K_{a,g}}L\alpha^a + \alpha^a \left\| \mathbf{D}\boldsymbol{\theta}^* \right\|_{a, \mathcal{G}}^a \right) \\ &\quad + \frac{\lambda\beta}{n} \sum_{l=1}^L \log \left\{ h \left(\left\| [\mathbf{D}\boldsymbol{\theta}^*]_{\mathcal{G}_l} \right\|_2 \right) \right\} + \frac{2K_{1,g}e^{3\sqrt{K_{a,g}}L\alpha^a} MC_{f,\ell}}{\mu}. \end{aligned}$$

According to (10.7), this completes the proof of Theorem 6.1. \square

10.2.2 Proof of Lemma 10.1

To prove Lemma 10.1, we need an intermediate result.

Lemma 10.2. *Let $s > L\sqrt{K_{a,g}}$. The following inequality holds*

$$\frac{1}{T} \int_{\left\{ \mathbf{u} \in \mathbb{R}^M : \left\| \mathbf{D}\mathbf{u} \right\|_{a, \mathcal{G}}^a > s \right\}} \prod_{l=1}^L g \left(\left\| [\mathbf{D}\mathbf{u}]_{\mathcal{G}_l} \right\|_2 \right) d\mathbf{u} \leq \frac{L^2 K_{a,g}}{(s - L\sqrt{K_{a,g}})^2},$$

where $T = \int_{\mathbb{R}^M} \prod_{l=1}^L g \left(\left\| [\mathbf{D}\mathbf{u}]_{\mathcal{G}_l} \right\|_2 \right) d\mathbf{u}$.

Proof of Lemma 10.2. Let \mathbf{U} be a random vector in \mathbb{R}^M with density

$$\mathbf{u} \mapsto \frac{1}{T} \prod_{l=1}^L g \left(\left\| [\mathbf{D}\mathbf{u}]_{\mathcal{G}_l} \right\|_2 \right),$$

where $T < \infty$ by Assumption 5.2-2. By Chebyshev inequality, we have

$$\begin{aligned} &\frac{1}{T} \int_{\left\{ \mathbf{u} \in \mathbb{R}^M : \left\| \mathbf{D}\mathbf{u} \right\|_{a, \mathcal{G}}^a > s \right\}} \prod_{l=1}^L g \left(\left\| [\mathbf{D}\mathbf{u}]_{\mathcal{G}_l} \right\|_2 \right) d\mathbf{u} = \mathbb{P} \left[\sum_{l=1}^L \left\| [\mathbf{D}\mathbf{U}]_{\mathcal{G}_l} \right\|_2^a > s \right] \\ &= \mathbb{P} \left[\sum_{l=1}^L \left(\left\| [\mathbf{D}\mathbf{U}]_{\mathcal{G}_l} \right\|_2^a - \mathbb{E} \left[\left\| [\mathbf{D}\mathbf{U}]_{\mathcal{G}_l} \right\|_2^a \right] \right) > s - \sum_{l=1}^L \mathbb{E} \left[\left\| [\mathbf{D}\mathbf{U}]_{\mathcal{G}_l} \right\|_2^a \right] \right] \\ &\leq \frac{\mathbb{E} \left[\left(\sum_{l=1}^L \left(\left\| [\mathbf{D}\mathbf{U}]_{\mathcal{G}_l} \right\|_2^a - \mathbb{E} \left[\left\| [\mathbf{D}\mathbf{U}]_{\mathcal{G}_l} \right\|_2^a \right] \right) \right)^2 \right]}{\left(s - \sum_{l=1}^L \mathbb{E} \left[\left\| [\mathbf{D}\mathbf{U}]_{\mathcal{G}_l} \right\|_2^a \right] \right)^2} \\ &= \frac{\text{var} \left(\sum_{l=1}^L \left\| [\mathbf{D}\mathbf{U}]_{\mathcal{G}_l} \right\|_2^a \right)}{\left(s - \sum_{l=1}^L \mathbb{E} \left[\left\| [\mathbf{D}\mathbf{U}]_{\mathcal{G}_l} \right\|_2^a \right] \right)^2} \\ &\leq \frac{\mathbb{E} \left[\left(\sum_{l=1}^L \left\| [\mathbf{D}\mathbf{U}]_{\mathcal{G}_l} \right\|_2^a \right)^2 \right]}{\left(s - \sum_{l=1}^L \mathbb{E} \left[\left\| [\mathbf{D}\mathbf{U}]_{\mathcal{G}_l} \right\|_2^a \right] \right)^2}. \end{aligned} \quad (10.22)$$

Next, by Cauchy-Schwartz inequality and Remark 5.1, we obtain

$$\mathbb{E} \left[\left(\sum_{l=1}^L \|[DU]_{\mathcal{G}_l}\|_2^a \right)^2 \right] \leq \mathbb{E} \left[L \sum_{l=1}^L \|[DU]_{\mathcal{G}_l}\|_2^{2a} \right] \leq L^2 K_{a,g} \quad (10.23)$$

and by Jensen inequality

$$s - \sum_{l=1}^L \mathbb{E} [\|[DU]_{\mathcal{G}_l}\|_2^a] \geq s - \sum_{l=1}^L \sqrt{\mathbb{E} [\|[DU]_{\mathcal{G}_l}\|_2^{2a}]} \geq s - L\sqrt{K_{a,g}} > 0. \quad (10.24)$$

Thus, combining (10.22), (10.23) and (10.24), we get

$$\frac{1}{T} \int_{\{\mathbf{u} \in \mathbb{R}^M : \|\mathbf{D}\mathbf{u}\|_{a,g}^a > s\}} \prod_{l=1}^L g(\|[DU]_{\mathcal{G}_l}\|_2) d\mathbf{u} \leq \frac{L^2 K_{a,g}}{(s - L\sqrt{K_{a,g}})^2}.$$

□

We now turn to the proof of Lemma 10.1

Proof of Lemma 10.1. Let us begin by the proof of inequality (10.20). We have

$$\begin{aligned} & \int_{\mathbb{R}^P} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 p_0(d\mathbf{u}) \\ &= \frac{1}{C_L \sqrt{\det(\mathbf{D}^T \mathbf{D})}} \int_{\mathbb{B}_{a,g}^{\alpha}(r_L) \cap \text{Im}(\mathbf{D})} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 \prod_{l=1}^L e^{-\alpha^a \|\mathbf{u}_{\mathcal{G}_l}\|_2^a} g(\|\mathbf{u}_{\mathcal{G}_l}\|_2) d\mathbf{u} \\ &\leq \frac{1}{C_L \sqrt{\det(\mathbf{D}^T \mathbf{D})}} \int_{\mathbb{B}_{a,g}^{\alpha}(r_L) \cap \text{Im}(\mathbf{D})} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 \prod_{l=1}^L g(\|\mathbf{u}_{\mathcal{G}_l}\|_2) d\mathbf{u}. \end{aligned} \quad (10.25)$$

In the following, we show inequality (10.20) by bounding the right term of inequality (10.25). By Lemma 4.1 and Remark 5.1, we get

$$\begin{aligned} & \frac{1}{\sqrt{\det(\mathbf{D}^T \mathbf{D})}} \frac{\int_{\mathbb{B}_{a,g}^{\alpha}(r_L) \cap \text{Im}(\mathbf{D})} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 \prod_{l=1}^L g(\|\mathbf{u}_{\mathcal{G}_l}\|_2) d\mathbf{u}}{\int_{\mathbb{R}^M} \prod_{l=1}^L g(\|[DU]_{\mathcal{G}_l}\|_2) d\mathbf{u}} \\ &\leq \frac{1}{\sqrt{\det(\mathbf{D}^T \mathbf{D})}} \frac{\int_{\text{Im}(\mathbf{D})} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 \prod_{l=1}^L g(\|\mathbf{u}_{\mathcal{G}_l}\|_2) d\mathbf{u}}{\int_{\mathbb{R}^M} \prod_{l=1}^L g(\|[DU]_{\mathcal{G}_l}\|_2) d\mathbf{u}} \\ &= \frac{\int_{\mathbb{R}^M} \|[DU]_{\mathcal{G}_1}\|_2^2 \prod_{l=1}^L g(\|[DU]_{\mathcal{G}_l}\|_2) d\mathbf{u}}{\int_{\mathbb{R}^M} \prod_{l=1}^L g(\|[DU]_{\mathcal{G}_l}\|_2) d\mathbf{u}} \leq K_{1,g}. \end{aligned}$$

Then

$$\frac{1}{\sqrt{\det(\mathbf{D}^T \mathbf{D})}} \int_{\mathbb{B}_{a,g}^{\alpha}(r_L) \cap \text{Im}(\mathbf{D})} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 \prod_{l=1}^L g(\|\mathbf{u}_{\mathcal{G}_l}\|_2) d\mathbf{u} \leq K_{1,g} T. \quad (10.26)$$

We now bound C_L^{-1} . By a change of variables, we obtain

$$\begin{aligned}
C_L^{-1} &= \left(\int_{\Theta_{p_0^D}} \prod_{l=1}^L e^{-\alpha^a \|\mathbf{D}\boldsymbol{\theta} - \mathbf{D}\boldsymbol{\theta}^*\|_{\mathcal{G}_l}} g(\|\mathbf{D}\boldsymbol{\theta} - \mathbf{D}\boldsymbol{\theta}^*\|_{\mathcal{G}_l}) d\boldsymbol{\theta} \right)^{-1} \\
&= \left(\int_{\{\mathbf{u} \in \mathbb{R}^M : \|\mathbf{D}\mathbf{u}\|_{a,\mathcal{G}}^a \leq r_L\}} e^{-\alpha^a \|\mathbf{D}\mathbf{u}\|_{a,\mathcal{G}}^a} \prod_{l=1}^L g(\|\mathbf{D}\mathbf{u}\|_{\mathcal{G}_l}) d\mathbf{u} \right)^{-1} \\
&\leq e^{r_L \alpha^a} \left(\int_{\{\mathbf{u} \in \mathbb{R}^M : \|\mathbf{D}\mathbf{u}\|_{a,\mathcal{G}}^a \leq r_L\}} \prod_{l=1}^L g(\|\mathbf{D}\mathbf{u}\|_{\mathcal{G}_l}) d\mathbf{u} \right)^{-1}.
\end{aligned}$$

Since $r_L = 3\sqrt{K_{a,g}}L > \sqrt{K_{a,g}}L$, Lemma 10.2 gives us

$$\begin{aligned}
C_L^{-1} &\leq e^{r_L \alpha^a} \left[T \left(1 - \frac{1}{T} \int_{\{\mathbf{u} \in \mathbb{R}^M : \|\mathbf{D}\mathbf{u}\|_{a,\mathcal{G}}^a > r_L\}} \prod_{l=1}^L g(\|\mathbf{D}\mathbf{u}\|_{\mathcal{G}_l}) d\mathbf{u} \right) \right]^{-1} \\
&\leq e^{r_L \alpha^a} T^{-1} \left(1 - \frac{L^2 K_{a,g}}{(r_L - L\sqrt{K_{a,g}})^2} \right)^{-1} \\
&= e^{r_L \alpha^a} T^{-1} \left(1 - \frac{1}{4} \right)^{-1} \leq 2e^{r_L \alpha^a} T^{-1}. \tag{10.27}
\end{aligned}$$

Combining (10.26) and (10.27), (10.25) becomes

$$\int_{\mathbb{R}^P} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 p_0(d\mathbf{u}) \leq 2K_{1,g} e^{r_L \alpha^a}.$$

That concludes the proof of inequality (10.20) in Lemma 10.1.

Next, we prove inequality (10.21). Remind that $\text{supp}(\pi) = \Theta$, $\text{supp}(p_0^D) = \Theta_{p_0^D}$. By (10.10), we get $\Theta_{p_0^D} \subseteq \Theta$ implying that p_0^D is absolutely continuous w.r.t. π . So $\text{KL}(p_0^D, \pi) < \infty$ which can be bounded. The bound in (10.21) can be proved as follows. By Lemma 4.1, we have

$$\begin{aligned}
\text{KL}(p_0^D, \pi) &= \int_{\mathbb{R}^M} \log \left(\frac{p_0^D(d\boldsymbol{\theta})}{\pi(d\boldsymbol{\theta})} \right) p_0^D(d\boldsymbol{\theta}) \\
&= \int_{\mathbb{R}^M} \log \left(\frac{C_{\alpha,g,R} \prod_{l=1}^L e^{-\alpha^a \|\mathbf{D}\boldsymbol{\theta} - \mathbf{D}\boldsymbol{\theta}^*\|_{\mathcal{G}_l}} g(\|\mathbf{D}\boldsymbol{\theta} - \mathbf{D}\boldsymbol{\theta}^*\|_{\mathcal{G}_l})}{C_L \prod_{l=1}^L e^{-\alpha^a \|\mathbf{D}\boldsymbol{\theta}\|_{\mathcal{G}_l}} g(\|\mathbf{D}\boldsymbol{\theta}\|_{\mathcal{G}_l})} \right) p_0^D(d\boldsymbol{\theta}) \\
&= \int_{\mathbb{R}^P} \log \left(\frac{C_{\alpha,g,R} \prod_{l=1}^L e^{\alpha^a \|\mathbf{t}_{\mathcal{G}_l}\|_2^a g(\|\mathbf{t}_{\mathcal{G}_l} - \mathbf{t}_{\mathcal{G}_l}^*\|_2)}{C_L \prod_{l=1}^L e^{\alpha^a \|\mathbf{t}_{\mathcal{G}_l} - \mathbf{t}_{\mathcal{G}_l}^*\|_2^a g(\|\mathbf{t}_{\mathcal{G}_l}\|_2)}} \right) p_0(dt) \\
&= \log \left(\frac{C_{\alpha,g,R}}{C_L} \right) + \alpha^a \sum_{l=1}^L \int_{\mathbb{R}^P} [\|\mathbf{t}_{\mathcal{G}_l}\|_2^a - \|\mathbf{t}_{\mathcal{G}_l} - \mathbf{t}_{\mathcal{G}_l}^*\|_2^a] p_0(dt) \\
&\quad + \sum_{l=1}^L \int_{\mathbb{R}^P} \log \left(\frac{g(\|\mathbf{t}_{\mathcal{G}_l} - \mathbf{t}_{\mathcal{G}_l}^*\|_2)}{g(\|\mathbf{t}_{\mathcal{G}_l}\|_2)} \right) p_0(dt),
\end{aligned}$$

where p_0 is a probability measure in \mathbb{R}^P defined in (10.12). We know that $\mathbf{t}^* = \mathbf{D}\boldsymbol{\theta}^*$, according to the fact that $\|\mathbf{t}_{\mathcal{G}_i}\|_2^a - \|\mathbf{t}_{\mathcal{G}_i} - \mathbf{t}_{\mathcal{G}_i}^*\|_2^a \leq \|\mathbf{t}_{\mathcal{G}_i}^*\|_2^a$ and Assumption 5.2-4, we get

$$\text{KL}(p_0^{\mathbf{D}}, \pi) \leq \log\left(\frac{C_{\alpha,g,R}}{C_L}\right) + \alpha^a \|\mathbf{D}\boldsymbol{\theta}^*\|_{a,g}^a + \lambda \sum_{l=1}^L \log\left\{h\left(\|[\mathbf{D}\boldsymbol{\theta}^*]_{\mathcal{G}_l}\|_2\right)\right\}. \quad (10.28)$$

Now, it remains to bound $\log(C_{\alpha,g,R}/C_L)$. Remind that $C_{\alpha,g,R}$ is the normalization factor of π , and thus we get

$$C_{\alpha,g,R} = \int_{\Theta} \prod_{l=1}^L \exp\left(-\alpha^a \|[\mathbf{D}\boldsymbol{\theta}]_{\mathcal{G}_l}\|_2^a\right) g\left(\|[\mathbf{D}\boldsymbol{\theta}]_{\mathcal{G}_l}\|_2\right) d\boldsymbol{\theta} \leq \int_{\mathbb{R}^M} \prod_{l=1}^L g\left(\|[\mathbf{D}\boldsymbol{\theta}]_{\mathcal{G}_l}\|_2\right) d\boldsymbol{\theta} = T.$$

Combining this with the bound of C_L^{-1} in (10.27), we obtain

$$\log\left(\frac{C_{\alpha,g,R}}{C_L}\right) \leq r_L \alpha^a + \log(2) \leq 1 + r_L \alpha^a. \quad (10.29)$$

Inserting (10.29) into (10.28), we get inequality (10.21). This completes the proof. \square

10.3 Proofs of Section 7

10.3.1 Proof of Lemma 7.1

In view of (H.1), V is also proper lsc and bounded from below. This together with (H.2) shows that V is prox-bounded for all $\gamma \in]0, \frac{\beta}{2\|\mathbf{X}\|^2}[$, and for any $\boldsymbol{\theta}$, $\frac{1}{2\gamma}\|\boldsymbol{\theta} - \cdot\|_{M_\gamma}^2 + V$ is proper lsc and level-bounded uniformly in $(\boldsymbol{\theta}, \gamma)$. This entails that the set of minimizers of this function, i.e. $\text{prox}_{\gamma V}^{M_\gamma}(\boldsymbol{\theta})$, is a non-empty compact set. The same reasoning leads to compactness of $\text{prox}_{\gamma J}(\boldsymbol{\theta})$.

Now, we have

$$\begin{aligned} \text{prox}_{\gamma V}^{M_\gamma}(\boldsymbol{\theta}) &= \underset{\mathbf{w} \in \mathbb{R}^M}{\text{Argmin}} \frac{1}{2\gamma} \|\mathbf{w} - \boldsymbol{\theta}\|_{M_\gamma}^2 + V(\mathbf{w}) \\ &= \underset{\mathbf{w} \in \mathbb{R}^M}{\text{Argmin}} \frac{1}{2} \|\mathbf{w} - \boldsymbol{\theta}\|_2^2 - \frac{\gamma}{\beta} \|\mathbf{X}(\mathbf{w} - \boldsymbol{\theta})\|_2^2 + \frac{\gamma}{\beta} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \gamma G(\mathbf{w}). \end{aligned}$$

By the Pythagoras relation, we then get

$$\begin{aligned} \text{prox}_{\gamma V}^{M_\gamma}(\boldsymbol{\theta}) &= \underset{\mathbf{w} \in \mathbb{R}^M}{\text{Argmin}} \frac{1}{2} \|\mathbf{w} - \boldsymbol{\theta}\|_2^2 + \frac{\gamma}{\beta} \left(\frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 - (\mathbf{X}(\boldsymbol{\theta} - \mathbf{w}))^T (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) \right) + \gamma G(\mathbf{w}) \\ &= \underset{\mathbf{w} \in \mathbb{R}^M}{\text{Argmin}} \frac{1}{2} \|\mathbf{w} - \boldsymbol{\theta}\|_2^2 - \frac{\gamma}{\beta} (\mathbf{w} - \boldsymbol{\theta})^T (\mathbf{X}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})) + \gamma G(\mathbf{w}) \\ &= \underset{\mathbf{w} \in \mathbb{R}^M}{\text{Argmin}} \frac{1}{2} \|\mathbf{w} - (\boldsymbol{\theta} - 2\gamma/\beta \mathbf{X}^T (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}))\|_2^2 + \gamma G(\mathbf{w}) \\ &= \text{prox}_{\gamma G}(\boldsymbol{\theta} - \gamma \nabla F(\boldsymbol{\theta})). \end{aligned}$$

For the last claim, suppose that $\boldsymbol{\theta} \in \text{Argmin}(V) \neq \emptyset$ and bounded but $\boldsymbol{\theta} \notin \text{prox}_{\gamma V}^{M_\gamma}(\boldsymbol{\theta})$. Thus, for any $\mathbf{p} \in \text{prox}_{\gamma V}^{M_\gamma}(\boldsymbol{\theta})$, we have $\mathbf{p} \neq \boldsymbol{\theta}$ and

$$V(\mathbf{p}) < \frac{1}{2\gamma} \|\mathbf{p} - \boldsymbol{\theta}\|_{M_\gamma}^2 + V(\mathbf{p}) \leq V(\boldsymbol{\theta}),$$

leading to a contradiction with $\boldsymbol{\theta}$ is a minimizer of V . \square

10.3.2 Proof of Lemma 7.2

The continuity and finiteness properties follow from [50, Theorem 1.17(c)] (see also [50, Theorem 1.25]), where we use (H.1) and (H.2). For the second claim, we have $\forall \boldsymbol{\theta} \in \mathbb{R}^M$

$$-\infty < \inf V \leq^{M_{\gamma, \gamma}} V(\boldsymbol{\theta}) \leq V(\boldsymbol{\theta}).$$

Moreover, let $\mathbf{p} \in \text{prox}_{\gamma V}^{M_{\gamma}}(\boldsymbol{\theta})$. Then, $\forall \delta] \gamma, \frac{\beta}{2\|\mathbf{X}\|^2} [$,

$$\begin{aligned} M_{\delta, \delta} V(\boldsymbol{\theta}) &= \inf_{\mathbf{w} \in \mathbb{R}^M} \frac{1}{2\delta} \|\mathbf{w} - \boldsymbol{\theta}\|_{M_{\delta}}^2 + V(\mathbf{w}) \\ &\leq \frac{1}{2\delta} \|\mathbf{p} - \boldsymbol{\theta}\|_{M_{\delta}}^2 + V(\mathbf{p}) \\ &\leq \frac{1}{2\gamma} \|\mathbf{p} - \boldsymbol{\theta}\|_{M_{\gamma}}^2 + V(\mathbf{p}) \\ &=^{M_{\gamma, \gamma}} V(\boldsymbol{\theta}). \end{aligned}$$

This together with continuity concludes the proof. \square

10.3.3 Proof of Lemma 7.3

By virtue of Lemma 7.1 and (H.3), $\text{prox}_{\gamma V}^{M_{\gamma}}$ is clearly non-empty and single valued. The continuity property follows from [50, Theorem 1.17(b)] (see also [50, Theorem 1.25]) and single-valuedness. By Lemma 7.2, V is prox-bounded with threshold $\beta/(2\|\mathbf{X}\|^2)$ and $^{M_{\gamma, \gamma}} V(\boldsymbol{\theta})$ is finite. Invoking [50, Example 10.32], we get that $-^{M_{\gamma, \gamma}} V$ is strictly continuous, subdifferentially regular and

$$\partial(-^{M_{\gamma, \gamma}} V)(\boldsymbol{\theta}) = \left\{ \gamma^{-1} M_{\gamma} \left(\text{prox}_{\gamma V}^{M_{\gamma}}(\boldsymbol{\theta}) - \boldsymbol{\theta} \right) \right\}, \quad \forall \boldsymbol{\theta} \in \mathbb{R}^M.$$

Combining this with [50, Theorem 9.18] applied to $-^{M_{\gamma, \gamma}} V$, we obtain that $^{M_{\gamma, \gamma}} V$ is differentiable and its gradient is precisely as given. \square

10.3.4 Proof of Lemma 7.4

- (i) Assumption (H.4), [50, Exercise 8.8(c)] and convexity of F imply that ∂V is hypomonotone of modulus r . In turn $S = \partial V + r' M_{\gamma} = \partial \left(V + \frac{r'}{2} \|\cdot\|_{M_{\gamma}}^2 \right)$ is monotone with $r' = \frac{r}{1-2\gamma\|\mathbf{X}\|^2/\beta}$, or equivalently that $V + \frac{r'}{2} \|\cdot\|_{M_{\gamma}}^2$ is convex. Let $\delta = \gamma^{-1} - r'$ and $W(\mathbf{w}, \boldsymbol{\theta}) = V(\mathbf{w}) + \frac{r'}{2} \|\mathbf{w} - \boldsymbol{\theta}\|_{M_{\gamma}}^2$. Thus

$$V(\mathbf{w}) + \frac{1}{2\gamma} \|\mathbf{w} - \boldsymbol{\theta}\|_{M_{\gamma}}^2 = W(\mathbf{w}, \boldsymbol{\theta}) + \frac{\delta}{2} \|\mathbf{w} - \boldsymbol{\theta}\|_{M_{\gamma}}^2.$$

$W(\cdot, \boldsymbol{\theta})$ is a convex function on \mathbb{R}^M and assumption (H.5) is equivalent to $\delta > 0$. Altogether, this entails that $W(\cdot, \boldsymbol{\theta}) + \frac{\delta}{2} \|\cdot - \boldsymbol{\theta}\|_{M_{\gamma}}^2$ is strongly convex uniformly in $\boldsymbol{\theta}$ and γ complying with (H.5). It then follows that $\text{prox}_{\gamma V}^{M_{\gamma}}$ is single-valued. We have

$$M_{\gamma} + \gamma \partial V = \gamma(\delta M_{\gamma} + S) = \gamma \delta (M_{\gamma} + \delta^{-1} S).$$

By Fermat's rule, we then get

$$\text{prox}_{\gamma V}^{M_{\gamma}} = (M_{\gamma} + \gamma \partial V)^{-1} \circ M_{\gamma} = (M_{\gamma} + \delta^{-1} S)^{-1} \circ (\gamma \delta)^{-1} M_{\gamma},$$

and the latter operator is well-defined as a single-valued operator since S is maximal monotone; see [3, Proposition 3.22 (ii)(d)]. Let $\mathbf{p} = \text{prox}_{\gamma V}^{M_\gamma}(\boldsymbol{\theta})$ and $\mathbf{p}' = \text{prox}_{\gamma V}^{M_\gamma}(\boldsymbol{\theta}')$. It then follows that

$$\mathbf{M}_\gamma \boldsymbol{\theta} - \gamma \delta \mathbf{M}_\gamma \mathbf{p} \in \gamma S(\mathbf{p}) \text{ and } \mathbf{M}_\gamma \boldsymbol{\theta}' - \gamma \delta \mathbf{M}_\gamma \mathbf{p}' \in \gamma S(\mathbf{p}'),$$

and monotonicity of S yields

$$(\mathbf{p}' - \mathbf{p})^T (\mathbf{M}_\gamma (\boldsymbol{\theta}' - \boldsymbol{\theta})) \geq \gamma \delta \|\mathbf{p}' - \mathbf{p}\|_{M_\gamma}^2 \geq \gamma \delta (1 - 2\gamma \|\mathbf{X}\|^2 / \beta) \|\mathbf{p}' - \mathbf{p}\|_2^2.$$

Using Cauchy-Schwartz inequality and $\|\mathbf{M}_\gamma\| \leq 1$, we then obtain

$$\|\mathbf{p}' - \mathbf{p}\|_2 \leq \kappa \|\boldsymbol{\theta}' - \boldsymbol{\theta}\|_2,$$

where $\kappa^{-1} = \gamma \delta (1 - 2\gamma \|\mathbf{X}\|^2 / \beta) = 1 - \gamma (r + 2\|\mathbf{X}\|^2 / \beta)$.

Let now $\mathbf{p} = \text{prox}_{\gamma V}^{M_\gamma}(\mathbf{M}_\gamma^{-1} \boldsymbol{\theta})$ and $\mathbf{p}' = \text{prox}_{\gamma V}^{M_\gamma}(\mathbf{M}_\gamma^{-1} \boldsymbol{\theta}')$. With the same arguments as before, we get

$$(\mathbf{p}' - \mathbf{p})^T (\boldsymbol{\theta}' - \boldsymbol{\theta}) \geq \gamma \delta \|\mathbf{p}' - \mathbf{p}\|_{M_\gamma}^2 \geq 0,$$

which shows monotonicity of $\text{prox}_{\gamma V}^{M_\gamma}$ (in fact even co-coercivity).

- (ii) Since $\text{prox}_{\gamma V}^{M_\gamma}$ is single-valued and Lipschitz continuous, C^1 smoothness and Lipschitz continuity of the gradient follow from Lemma 7.3. Let's now turn to hypomonotonicity. We have

$$\begin{aligned} \mathbf{p} &= \text{prox}_{\gamma V}^{M_\gamma}(\boldsymbol{\theta}) = (\mathbf{M}_\gamma + \delta^{-1} S)^{-1} \left((\gamma \delta)^{-1} \mathbf{M}_\gamma \boldsymbol{\theta} \right) \\ &\iff \gamma^{-1} \mathbf{M}_\gamma \boldsymbol{\theta} - \delta \mathbf{M}_\gamma \mathbf{p} \in S(\mathbf{p}) \\ &\iff \mathbf{p} \in S^{-1}(\gamma^{-1} \mathbf{M}_\gamma \boldsymbol{\theta} - \delta \mathbf{M}_\gamma \mathbf{p}) \\ &\iff (\delta \gamma)^{-1} \mathbf{M}_\gamma \boldsymbol{\theta} \in \mathbf{M}_\gamma \left((\delta \mathbf{M}_\gamma)^{-1} + S^{-1} \right) (\gamma^{-1} \mathbf{M}_\gamma \boldsymbol{\theta} - \delta \mathbf{M}_\gamma \mathbf{p}). \end{aligned}$$

Since S is monotone, so is S^{-1} . In turn, with similar arguments as for S above, $\left((\delta \mathbf{M}_\gamma)^{-1} + S^{-1} \right)^{-1}$ is well-defined as a single-valued operator, and is monotone. Therefore,

$$\gamma^{-1} \mathbf{M}_\gamma \boldsymbol{\theta} - \delta \mathbf{M}_\gamma \mathbf{p} = \left((\delta \mathbf{M}_\gamma)^{-1} + S^{-1} \right)^{-1} \left((\delta \gamma)^{-1} \boldsymbol{\theta} \right).$$

Using Lemma 7.3, we obtain

$$\delta \gamma \nabla^{M_\gamma} V(\boldsymbol{\theta}) = (\delta - \gamma^{-1}) \mathbf{M}_\gamma \boldsymbol{\theta} + \left((\delta \mathbf{M}_\gamma)^{-1} + S^{-1} \right)^{-1} \left((\delta \gamma)^{-1} \boldsymbol{\theta} \right).$$

As we argued above, on the one hand, $\left((\delta \mathbf{M}_\gamma)^{-1} + S^{-1} \right)^{-1}$ is monotone. On the other hand, $\delta - \gamma^{-1} = -r'$. Clearly $(\delta - \gamma^{-1}) \mathbf{M}_\gamma$ is symmetric definite negative, whence we deduce the hypomonotonicity claim. \square

10.3.5 Proof of Proposition 7.1

Recall that we denote with the same symbol the measure and its density with respect to the Lebesgue measure. Thus

$$\|\mu_\gamma - \mu_n\|_{\text{TV}} = \int_{\mathbb{R}^M} |\mu_\gamma(\boldsymbol{\theta}) - \mu_n \boldsymbol{\theta}| d\boldsymbol{\theta},$$

where

$$\mu_\gamma(\boldsymbol{\theta}) = \exp\left(-\mathbf{M}_\gamma \gamma V(\mathbf{M}_\gamma^{-1/2} \boldsymbol{\theta})\right) / Z_\gamma \text{ and } \mu_n(\boldsymbol{\theta}) = \exp(-V(\boldsymbol{\theta})) / Z,$$

and $Z = \int_{\mathbb{R}^M} \exp(-V(\mathbf{u})) d\mathbf{u}$. We have $\det(\mathbf{M}_\gamma) \rightarrow 1$ as $\gamma \rightarrow 0$. In view of Lemma 7.2, applying the monotone convergence theorem, we conclude that $Z_\gamma \rightarrow Z$. This together with Lemma 7.2 yield that μ_γ converges to μ_n pointwise. We conclude using Scheffé(-Riesz) theorem [37, 53]. \square

10.3.6 Proof of Theorem 7.1

By Lemma 7.4(i) and Lemma 7.1, we have Lipschitz continuity of the operator $S = \frac{\mathbf{I}_M - \text{prox}_{\gamma G} \circ (\mathbf{I}_M - \gamma \nabla F)}{\gamma}$. Let κ be this Lipschitz constant (which depends on γ). By these two lemmas again, we have for any bounded $\boldsymbol{\theta} \in \text{Argmin}(V)$

$$\|S(\mathbf{w})\|_2^2 = \|S(\mathbf{w}) - S(\boldsymbol{\theta})\|_2^2 \leq \kappa^2 \|\mathbf{w} - \boldsymbol{\theta}\|_2^2 \leq 2\kappa^2 \|\mathbf{w}\|_2^2 + 2\kappa^2 \|\boldsymbol{\theta}\|_2^2 \leq c(1 + \|\mathbf{w}\|_2^2),$$

where $c = 2\kappa^2 \max(1, \|\boldsymbol{\theta}\|_2^2)$. Thus the global Lipschitz and growth conditions required on the drift coefficient are verified, and the bound follows from [56, Theorem 7.3, Chapter II] or [35, Theorem 10.2.2 and Remark 10.2.3]. Using Jensen's inequality we get the bound on the moment. \square

10.3.7 Proof of Lemma 7.5

This is a probably known result, for which we provide a simple proof. Since $J = G \circ \mathbf{D}^{-1}$ is separable and w is continuous and lower-bounded, we have

$$\min_{\mathbf{w} \in \mathbb{R}^P} \frac{1}{2} \|\mathbf{w} - \mathbf{u}\|_2^2 + \gamma J(\mathbf{w}) = \sum_{l=1}^L \min_{\mathbf{v} \in \mathbb{R}^G} \frac{1}{2} \|\mathbf{v} - \mathbf{u}_{G_l}\|_2^2 + \gamma w(\|\mathbf{v}\|_2),$$

and thus, $\forall l \in \{1, \dots, L\}$,

$$[\text{prox}_{\gamma J}(\mathbf{u})]_{G_l} = \underset{\mathbf{v} \in \mathbb{R}^G}{\text{Argmin}} \frac{1}{2} \|\mathbf{v} - \mathbf{u}_{G_l}\|_2^2 + \gamma w(\|\mathbf{v}\|_2). \quad (10.30)$$

If $\mathbf{u}_{G_l} = 0$, then as w is an increasing function, $[\text{prox}_{\gamma J}(\mathbf{u})]_{G_l} = 0$. For $\mathbf{u}_{G_l} \neq 0$, by isotropy of problem (10.30), we can write

$$\min_{\mathbf{v} \in \mathbb{R}^G} \frac{1}{2} \|\mathbf{v} - \mathbf{u}_{G_l}\|_2^2 + \gamma w(\|\mathbf{v}\|_2) = \min_{t \geq 0} \gamma w(t) + \left(\min_{\|\mathbf{v}\|_2=t} \frac{1}{2} \|\mathbf{v} - \mathbf{u}_{G_l}\|_2^2 \right). \quad (10.31)$$

The inner minimization problem amounts to solving for the orthogonal projector on the ℓ_2 sphere in \mathbb{R}^G of radius t , which is obviously $\mathbf{v} = t \frac{\mathbf{u}_{G_l}}{\|\mathbf{u}_{G_l}\|_2}$ since $\mathbf{u}_{G_l} \neq 0$. Inserting this into (10.31) and rearranging the terms, (10.30) becomes

$$[\text{prox}_{\gamma J}(\mathbf{u})]_{G_l} = \frac{\mathbf{u}_{G_l}}{\|\mathbf{u}_{G_l}\|_2} \underset{t \geq 0}{\text{Argmin}} \frac{1}{2} (t - \|\mathbf{u}_{G_l}\|_2)^2 + \gamma w(t) = \frac{\mathbf{u}_{G_l}}{\|\mathbf{u}_{G_l}\|_2} \text{prox}_{\gamma w}(\|\mathbf{u}_{G_l}\|_2),$$

where we used even-symmetry of w . \square

10.3.8 Proof of Lemma 7.6

We can see directly that

$$\tilde{w}(t) = \alpha^a t^a + c \log(\tau^b + t^b) - c \log(\tau^b)$$

is non negative, differentiable and non decreasing on $]0, +\infty[$. Let $u : t \in]0, +\infty[\mapsto t + \tilde{w}(t)$, it suffices to prove that this function is strictly unimodal to apply [2, Theorem 1]. Indeed, one can show by calculating the first and second derivatives of u that

$$\lim_{t \rightarrow 0^+} u'(t) = -\infty, \quad \lim_{t \rightarrow +\infty} u'(t) = 1,$$

and

$$u''(t) > 0, \quad \forall t \in]0, +\infty[.$$

These two arguments show that there exists a unique minimum of u on $]0, +\infty[$. So, we can apply [2, Theorem 1] to get the proximity of \tilde{w} which is the same as that of w since the minimization problem is unchanged by adding a constant. \square

Acknowledgement. This work was supported by Conseil Régional de Basse-Normandie and partly by Institut Universitaire de France.

References

- [1] Y. Amit and D. Geman. Shape quantization and recognition with randomized trees. *Neural Comput.*, 9(7):1545–1588, Oct. 1997.
- [2] A. Antoniadis and J. Fan. Regularization of Wavelet Approximations. *Journal of the American Statistical Association*, 96:939–967, 2001.
- [3] H. H. Bauschke, J. M. Borwein, and P. L. Combettes. Bregman monotone optimization algorithms. *SIAM Journal on Control and Optimization*, 42(2):596–636, 2003.
- [4] H. H. Bauschke and P. L. Combettes. *Convex analysis and monotone operator theory in Hilbert spaces*. Springer, 2011.
- [5] G. Biau. Analysis of a random forests model. *J. Mach. Learn. Res.*, 13(1):1063–1095, Apr. 2012.
- [6] G. Biau and L. Devroye. On the layered nearest neighbour estimate, the bagged nearest neighbour estimate and the random forest method in regression and classification. *J. Multivar. Anal.*, 101(10):2499–2518, Nov. 2010.
- [7] G. Biau, L. Devroye, and G. Lugosi. Consistency of random forests and other averaging classifiers. *J. Mach. Learn. Res.*, 9:2015–2033, June 2008.
- [8] L. Breiman. Bagging predictors. *Mach. Learn.*, 24(2):123–140, Aug. 1996.
- [9] L. Breiman. Random forests. *Mach. Learn.*, 45(1):5–32, Oct. 2001.
- [10] P. Bühlmann and S. Van De Geer. *Statistics for high-dimensional data: methods, theory and applications*. Springer, 2011.

- [11] E. Candès and D. Donoho. New tight frames of curvelets and optimal representations of objects with piecewise C^2 singularities. *Communications on Pure and Applied Mathematics*, 57(2):219–266, 2004.
- [12] E. J. Candès. Ridgelets: Estimating with Ridge Functions. *Annals of Statistics*, 31, 1999. 1561–1599.
- [13] L. Chaari, J.-Y. Tourneret, C. Chaux, and H. Batatia. A hamiltonian monte carlo method for non-smooth energy sampling. Technical Report arXiv:1401.3988, , 2014.
- [14] R. Coifman and D. Donoho. Translation invariant de-noising. In A. Antoniadis and G. Oppenheim, editors, *Wavelets and Statistics*. Springer-Verlag, 1995. 125–150.
- [15] A. Dalalyan and A. Tsybakov. Pac-bayesian bounds for the expected error of aggregation by exponential weights. Technical report, Université Paris 6, CREST and CERTIS, Ecole des Ponts ParisTech, 2009. personal communication.
- [16] A. Dalalyan and A. B. Tsybakov. Aggregation by exponential weighting, sharp pac-bayesian bounds and sparsity. *Mach. Learn.*, 72(1-2):39–61, Aug. 2008.
- [17] A. S. Dalalyan and J. Salmon. Sharp oracle inequalities for aggregation of affine estimators. *Ann. Statist.*, 40(4):2327–2355, 08 2012.
- [18] A. S. Dalalyan and A. B. Tsybakov. Aggregation by exponential weighting and sharp oracle inequalities. In *Proceedings of the 20th Annual Conference on Learning Theory*, COLT’07, pages 97–111, Berlin, Heidelberg, 2007. Springer-Verlag.
- [19] A. S. Dalalyan and A. B. Tsybakov. Mirror averaging with sparsity priors. *Bernoulli*, 18(3):914–944, 08 2012.
- [20] A. S. Dalalyan and A. B. Tsybakov. Sparse regression learning by aggregation and langevin monte-carlo. *J. Comput. Syst. Sci.*, 78(5):1423–1443, Sept. 2012.
- [21] D. L. Donoho. Observed universality of phase transitions in high-dimensional geometry, with implications for modern data analysis and signal processing. *Phil. Trans. Royal Soc. A*, 367(1906):4273–4293, 2009.
- [22] D. L. Donoho and I. M. Johnstone. Ideal spatial adaptation by wavelet shrinkage. *Biometrika*, 81:425–455, 1994.
- [23] D. L. Donoho and I. M. Johnstone. Adapting to unknown smoothness via wavelet shrinkage. *Journal of the American Statistical Association*, pages 1200–1224, 1995.
- [24] D. L. Donoho, I. M. Johnstone, G. Kerkyacharian, and D. Picard. Wavelet shrinkage: asymptopia. *Journal of the Royal Statistical Society, Ser. B*, pages 371–394, 1995.
- [25] A. Durmus and E. Moulines. Non-asymptotic convergence analysis for the Unadjusted Langevin Algorithm. Preprint hal-01176132, July 2015.
- [26] A. Durmus, E. Moulines, and M. Pereyra. Sampling from convex non continuously differentiable functions, when Moreau meets Langevin. Preprint hal-01267115, Feb. 2016.
- [27] K. Fang, S. Kotz, and K. Ng. *Symmetric multivariate and related distributions*. Monographs on statistics and applied probability. Chapman and Hall, 1990.

- [28] Y. Freund. Boosting a weak learning algorithm by majority. *Information and Computation*, 121(2):256–285, 1995.
- [29] R. Genuer. *Random Forests: elements of theory, variable selection and applications*. Theses, Université Paris Sud - Paris XI, Nov. 2010.
- [30] I. S. Gradshteyn and I. M. Ryzhik. *Table of Integrals, Series, and Products*. Academic Press, London, 4th edition, 1965.
- [31] D. Higham, X. Mao, and A. Stuart. Strong convergence of euler-type methods for nonlinear stochastic differential equations. *SIAM J. Numer. Anal.*, 40(3):1041–1063, 2003.
- [32] J. Huang, J. L. Horowitz, and F. Wei. Variable selection in nonparametric additive models. *Ann. Statist.*, 38(4):2282–2313, 08 2010.
- [33] N. Ikeda and S. Watanabe. *Stochastic differential equations and diffusion processes*. NH, 1989.
- [34] J. Kivinen and M. K. Warmuth. Exponentiated gradient versus gradient descent for linear predictors. *Inf. Comput.*, 132(1):1–63, Jan. 1997.
- [35] P. E. Kloeden and E. Platen. *Numerical solution of stochastic differential equations*. Stochastic Modelling and Applied Probability. Springer, 1995.
- [36] V. Koltchinskii and M. Yuan. Sparse recovery in large ensembles of kernel machines, 2008.
- [37] N. Kusolitsch. Why the theorem of scheffé should be rather called a theorem of riesz. *Periodica Mathematica Hungarica*, 61(1):225–229, 2010.
- [38] G. Lecué. Simultaneous adaptation to the margin and to complexity in classification. *Ann. Statist.*, 35(4):1698–1721, 08 2007.
- [39] G. Leung and A. R. Barron. Information theory and mixing least-squares regressions. *IEEE Trans. Inform. Theory*, 52(8):3396–3410, 2006.
- [40] N. Littlestone and M. K. Warmuth. The weighted majority algorithm. *Inf. Comput.*, 108(2):212–261, Feb. 1994.
- [41] J. Maly. Lectures on change of variables in integral. Preprint 305, Department of Mathematics, University of Helsinki, 2001.
- [42] L. Meier, S. V. D. Geer, P. Bühlmann, and E. T. H. Zivich. The group lasso for logistic regression. *Journal of the Royal Statistical Society, Series B*, 2008.
- [43] L. Meier, S. van de Geer, and P. Bühlmann. High-dimensional additive modeling. *Ann. Statist.*, 37(6B):3779–3821, 12 2009.
- [44] A. Nemirovski. Topics in non-parametric statistics, 2000.
- [45] M. Pereyra. Proximal markov chain monte carlo algorithms. *Statistics and Computing*, 26(4):745–760, 2016.

- [46] P. Ravikumar, J. Lafferty, H. Liu, and L. Wasserman. Sparse additive models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(5):1009–1030, 2009.
- [47] P. Rigollet and A. Tsybakov. Linear and convex aggregation of density estimators. *Mathematical Methods of Statistics*, 16(3):260–280, 2007.
- [48] P. Rigollet and A. B. Tsybakov. Sparse estimation by exponential weighting. *Statist. Sci.*, 27(4):558–575, 11 2012.
- [49] G. O. Roberts and R. L. Tweedie. Exponential Convergence of Langevin Distributions and Their Discrete Approximations. *Bernoulli*, 2(4):341–363, 1996.
- [50] R. T. Rockafellar and R. Wets. *Variational analysis*, volume 317. Springer Verlag, 1998.
- [51] L. I. Rudin, S. Osher, and E. Fatemi. Nonlinear total variation based noise removal algorithms. *Phys. D*, 60(1-4):259–268, Nov. 1992.
- [52] R. E. Schapire. The strength of weak learnability. *Mach. Learn.*, 5(2):197–227, July 1990.
- [53] H. Scheffé. A useful convergence theorem for probability distributions. *Ann. Math. Statist.*, 18(3):434–438, 09 1947.
- [54] R. Tibshirani, M. Saunders, S. Rosset, J. Zhu, and K. Knight. Sparsity and smoothness via the fused lasso. *Journal of the Royal Statistical Society Series B*, pages 91–108, 2005.
- [55] V. G. Vovk. Aggregating strategies. In *Proceedings of the Third Annual Workshop on Computational Learning Theory, COLT '90*, pages 371–386, San Francisco, CA, USA, 1990. Morgan Kaufmann Publishers Inc.
- [56] M. Xuerong. *Stochastic differential equations and applications*. Woodhead Publishing, 2007.
- [57] Y. Yang. Aggregating regression procedures to improve performance. *Bernoulli*, 10(1):25–47, 2004.
- [58] M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society, Series B*, 68:49–67, 2006.