



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

Duy Tung Luu, Jalal M. Fadili, Christophe Chesneau

► To cite this version:

Duy Tung Luu, Jalal M. Fadili, Christophe Chesneau. PAC-Bayesian risk bounds for group-analysis sparse regression by exponential weighting. Journal of Multivariate Analysis, In press. hal-01367742v4

HAL Id: hal-01367742

<https://hal.science/hal-01367742v4>

Submitted on 22 May 2018

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

Normandie Univ, ENSICAEN, CNRS, GREYC, France

Christophe Chesneau

Normandie Univ, UNICAEN, CNRS, LMNO, France

Abstract

In this paper, we consider a high-dimensional non-parametric regression model with fixed design and i.i.d. random errors. We propose an 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 convergence guarantees. In turn, this allows us to implement our estimator and validate it on some numerical experiments.

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

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) are random errors. (1.1) is equivalently written in vector form $\mathbf{Y} = \mathbf{f} + \boldsymbol{\xi}$. Assume that there exists a dictionary $\mathcal{H} = \{f_j : \mathcal{X} \rightarrow \mathbb{R}, j \in \{1, \dots, p\}\}$ such that f is well approximated by a linear combination of elements in \mathcal{H} , i.e., there exists $\bar{\boldsymbol{\theta}} = (\bar{\theta}_1, \dots, \bar{\theta}_p)^\top \in \mathbb{R}^p$ such that $f_{\bar{\boldsymbol{\theta}}} = \sum_{j=1}^p \bar{\theta}_j f_j$ 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., $p \gg n$). Thus, in such a setting, the classical least-squares to estimate $\bar{\boldsymbol{\theta}}$ is obviously not applicable.

The idea of aggregating elements in a dictionary has been introduced in machine learning to combine different techniques (see [39, 56]) with some procedures such as bagging [7], boosting [31, 54] and random forests [1, 4–6, 8, 32]. In the recent years, there has been a flurry of research on the use of the concept of sparsity in various areas

Email addresses: `duy-tung.luu@ensicaen.fr` (Tung Duy Luu), `Jalal.Fadili@greyc.ensicaen.fr` (Jalal Fadili), `christophe.chesneau@unicaen.fr` (Christophe Chesneau)

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 to be much smaller than n . This makes it possible to build an estimate $\hat{f}_{\tilde{\theta}}$ 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 [25–27] on orthogonal wavelet thresholding estimators. Oracle inequalities (according to the terminology introduced in, e.g., [26]), which are at the heart of our work, quantify the quality of an estimator compared to the best possible one among a collection of estimators. 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 performance of the best model of aggregation in a given class Θ (in the probabilistic sense). This idea is expressed in the following type of inequalities:

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

where $C \geq 1$ and the remainder term $\Delta_{n,p}(\theta)$ depends on the performance of the estimator, the complexity of θ , the dimension p and the sample size n . Such type of inequality is called balanced oracle inequality. Inequalities of type (1.2) are well adapted under the sparsity scenario. Namely, the complexity of θ in the remaining term is characterized by the sparsity parameters (like the number of its non-zero components), in which case these inequalities 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”), $\Delta_{n,p}(\theta)$ is small even if $n \ll p$ 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 $Y = y$, The penalization approach considers the minimization problem

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

where $\text{pen} : \mathbb{R}^p \rightarrow \mathbb{R}^+$ is a sparsity promoting penalty function, see, e.g., [9]. Our work focuses on EWA approach that we briefly describe now.

1.2.2. Exponential Weighted Aggregation (EWA)

Let (Λ, \mathcal{A}) be 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

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

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

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

This idea was initially proposed in [35, 39, 56] with a uniform prior on a finite set Λ . Observe that (1.4) can be interpreted as the Bayesian posterior conditional mean in the regression model but only when the noise is Gaussian noise $\mathcal{N}(0, \sigma^2 \mathbf{I}_n)$ and provided that $\beta = 2\sigma^2$ and the prior density is π .

In the literature, the works in [18, 19, 21, 23] consider deterministic dictionaries. These papers proposed several PAC-Bayesian¹ type of oracle inequalities under different assumptions. Especially, the assumptions in [23] 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 $O(\|\theta\|_0 \ln(p)/n)$, which scales linearly with the sparsity level and increases in p only logarithmically.

The random dictionary case are tackled in [45]. The initial idea is to obtain two independent samples from the initial sample by randomization or sample splitting (see [37, 49, 58]). 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 [38]. More recent work such as [20] generalizes the idea where the pre-estimators are affine and the priors are continuous.

A shortcoming of EWA is its suboptimality in deviation. In particular, the work in [16, Section 2] has shown that the EWA leads a suboptimal remainder term for oracle inequalities in probability.

1.2.3. Generalization of sparsity assumption

Analysis sparsity Let $q \geq p$ and $\mathbf{D}^T \in \mathbb{R}^{q \times p}$ be a linear analysis operator. The analysis sparsity assumption means that $\|\mathbf{D}^T \theta\|_0 \ll n$. A typical example is total variation [53] where the operator \mathbf{D}^T corresponds to “finite differences” (i.e., $(\mathbf{D}^T \theta)_1 = \theta_1$, $(\mathbf{D}^T \theta)_j = \theta_j - \theta_{j-1}$, $\forall j \geq 2$). Another example is the fused Lasso [55] where \mathbf{D}^T is a positive combination of the identity and finite differences.

Group sparsity Group sparsity corresponds to saying that the aggregator θ is block sparse, see Section 4.1 for formal details. Group sparsity is at the heart of the group Lasso and related methods [34, 36, 42, 43, 48, 59]. In the EWA context, the group sparsity prior is considered in [50] 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 a frame, and thus is not necessarily invertible unlike previous work. In a finite space, that equivalents to the fact that $\mathbf{D}\mathbf{D}^T$ is invertible. In addition, our prior class (see (4.2)) is much more general than previously proposed ones [23] 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 Assumptions (G.1)-(G.4)).
- We establish a group-analysis SOI where the remainder term depends on the number of active groups in $\mathbf{D}^T \theta$ and on the function g (see Theorem 5.1)
- 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 $O(\|\mathbf{D}^T \theta\|_{0,\mathcal{G}} \ln(L)/n)$, where $\|\mathbf{D}^T \theta\|_{0,\mathcal{G}}$ is the number of active groups in $\mathbf{D}^T \theta$, and L is the total number of groups (see Corollary 5.2). This rate coincides with the classical one $O(\|\theta\|_0 \ln(p)/n)$ under the sparsity scenario, i.e., $\mathbf{D} = \mathbf{I}_p$ and $L = p$.

¹PAC stands for Probably Approximately Correct. A PAC-bound is probably correct as it not a deterministic guarantee allowing a small probability that the estimator does not behave well. It is also approximately correct as it is tolerant to an inexact performance of the estimator.

- We also propose a forward-backward proximal Langevin Monte-Carlo (LMC) algorithm to sample from the target distribution (6.3) (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 [23] which are a classical starting point in literature for EWA in the deterministic case. In Section 4, we describe our EWA procedure after specifying the aggregation dictionary and our prior family. In Section 5, we establish our main results, namely group-analysis SOI. Section 6 is devoted to the forward-backward proximal LMC algorithm that implements EWA, and the numerical experiments on several numerical settings are described in Section 7. The proofs of all results are collected in Section 9.

2. Notations and Preliminaries

Before proceeding, let us introduce some notations and definitions.

Vectors and matrices For a d -dimensional Euclidean space \mathbb{R}^d , we endow it with its usual inner product $\langle \cdot, \cdot \rangle$ and associated norm $\|\cdot\|_2$. \mathbf{I}_d is the identity matrix on \mathbb{R}^d . For $q > 0$ and $\mathbf{x} \in \mathbb{R}^d$, we also define $\|\mathbf{x}\|_q = \left(\sum_{j=1}^d |x_j|^q\right)^{1/q}$ with the usual adaptation $\|\mathbf{x}\|_\infty = \max_{j \in \{1, \dots, d\}} |x_j|$. It is the ℓ_q norm for $q \geq 1$, and quasi-norm for $q \in]0, 1[$. $\|\mathbf{x}\|_0$ is the ℓ_0 pseudo-norm which counts the number of non-zero elements in \mathbf{x} . Let $\mathbf{M} \in \mathbb{R}^{d \times d}$ symmetric positive definite, we denote $\langle \cdot, \cdot \rangle_{\mathbf{M}} = \langle \cdot, \mathbf{M} \cdot \rangle$ and $\|\cdot\|_{\mathbf{M}}$ its associated norm. Of course, $\|\cdot\|_{\mathbf{M}}$ and $\|\cdot\|_2$ are equivalent.

For a matrix $\mathbf{M} \in \mathbb{R}^{d \times r}$, we set $\sigma(\mathbf{M}) = (\sigma_1(\mathbf{M}), \dots, \sigma_r(\mathbf{M}))^\top \in \mathbb{R}^r$ be the vector of singular values of \mathbf{M} in non-increasing order. Note that, when \mathbf{M} is symmetric semi-definite positive, $\sigma(\mathbf{M})$ is also the ordered vector of positive eigenvalues of \mathbf{M} . We denote $\|\mathbf{M}\|$ the spectral norm of \mathbf{M} .

For $\mathcal{I} \subset \{1, \dots, d\}$, $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{M} \in \mathbb{R}^{d \times d}$, $\mathbf{x}_{\mathcal{I}}$ is the subvector whose entries are those of \mathbf{x} restricted to the indices in \mathcal{I} , and $\mathbf{M}_{\mathcal{I}}$ the submatrix whose rows and columns are those of \mathbf{M} indexed by \mathcal{I} .

Let us denote $\text{vec} : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{d^2}$ the vectorization operator of a matrix. For any matrix \mathbf{M} , \mathbf{M}^\top denotes its transpose. For any square matrix \mathbf{M} , $\det(\mathbf{M})$ is its determinant.

Definition 2.1 (Frame). A matrix $\mathbf{M} \in \mathbb{R}^{d \times r}$ is a frame if there exist two constants ν and κ with $\nu \geq \kappa > 0$, called frame bounds, such that the generalized Parseval relation is satisfied, i.e., $\kappa \|\mathbf{x}\|_2^2 \leq \|\mathbf{M}^\top \mathbf{x}\|_2^2 \leq \nu \|\mathbf{x}\|_2^2$, $\forall \mathbf{x} \in \mathbb{R}^d$.

By the Courant-Fischer theorem, Definition 2.1 is equivalent to the fact that κ (resp. ν) is a lower (resp. upper) bound of the eigenvalues of $\mathbf{M}\mathbf{M}^\top$. Moreover, since $\kappa > 0$, we have that $\mathbf{M}\mathbf{M}^\top$ is bijective and \mathbf{M}^\top is injective. The frame is said tight when $\kappa = \nu$. Typical examples of (tight) frames that have been used in statistics are translation invariant wavelets [15], ridgelets [11] and curvelets [10] (example of groups and what they represent for wavelets/ridgelets/curvelets in applications are discussed in [14]). Let $\tilde{\mathbf{M}} \in \mathbb{R}^{d \times r}$ be the canonical dual frame associated to \mathbf{M} , i.e., $\tilde{\mathbf{M}} = (\mathbf{M}\mathbf{M}^\top)^{-1} \mathbf{M}$. We know that

$$\tilde{\mathbf{M}}\mathbf{M}^\top = \mathbf{I}_d \quad (2.1)$$

and

$$\frac{1}{\kappa} \geq \sigma_1(\tilde{\mathbf{M}}^\top \tilde{\mathbf{M}}) \geq \dots \geq \sigma_d(\tilde{\mathbf{M}}^\top \tilde{\mathbf{M}}) \geq \frac{1}{\nu}. \quad (2.2)$$

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.

Functions For a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, we define $\mathbf{f} = (f(x_1), \dots, f(x_d))^T$ its d -sample vector form, $\|\mathbf{f}\|_d = \left(\frac{1}{d} \sum_{j=1}^d f^2(x_j)\right)^{1/2}$ and $\|f\|_\infty = \sup_{\mathbf{x} \in \mathbb{R}^d} |f(\mathbf{x})|$.

For a function $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$, its effective domain is $\text{dom}(f) = \{\mathbf{x} \in \mathbb{R}^d : f(\mathbf{x}) < +\infty\}$ and f is proper if $f(\mathbf{x}) > -\infty$ for all \mathbf{x} and $\text{dom}(f) \neq \emptyset$ as is the case when it is finite-valued. A function $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{-\infty, +\infty\}$ is lower semi continuous (lsc) at \mathbf{x}_0 if $\liminf_{\mathbf{x} \rightarrow \mathbf{x}_0} f(\mathbf{x}) \geq f(\mathbf{x}_0)$. For a differentiable function f , ∇f is its (Euclidean) gradient.

For a set Ω , I_Ω is its characteristic function, i.e., 1 if the argument is in Ω and 0 otherwise. Define $\text{sgn} : \mathbb{R} \mapsto \{-1, 1\}$ be the sign function, i.e., $\text{sgn}(x) = 1$ when $x \geq 0$, and $\text{sgn}(x) = -1$ otherwise. Let $x \in \mathbb{R}$, define $\lfloor x \rfloor$ be the stands of integer part of x . Recall the Gamma function $\Gamma :]0, +\infty[\rightarrow]0, +\infty[$, $\Gamma : t \mapsto \int_0^\infty x^{t-1} \exp(-x) dx$.

Definition 2.2 (Proximal mapping and Moreau envelope). *Let $\mathbf{M} \in \mathbb{R}^{d \times d}$ be symmetric positive definite. For a proper lsc function f and $\gamma > 0$, the proximal mapping and Moreau envelope in the metric \mathbf{M} are defined respectively by*

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

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

Useful integration formulas The following lemmas contain useful formula used throughout the paper.

Lemma 2.1 ([33, 3.251.11]). *Let $p, \gamma, \nu, \eta > 0$. If $\gamma/\nu < \eta + 1$ we have*

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

otherwise this integral is not definite.

Lemma 2.2 (Cartesian to spherical coordinates [30]). *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 then have*

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

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

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 5.1).

Lemma 2.3. *Let $\Theta \subseteq \mathbb{R}^d$ be a measurable set. For any $\mathbf{M} \in \mathbb{R}^{r \times d}$ corresponding to the analysis operator of a frame of \mathbb{R}^d , let $u : \mathbb{R}^r \rightarrow \mathbb{R}$ such that the mapping $\mathbf{x} \mapsto u(\mathbf{M}\mathbf{x})$ is measurable on Θ . We have*

$$\int_{\Theta} u(\mathbf{M}\mathbf{x}) d\mathbf{x} = \frac{1}{\sqrt{\det(\mathbf{M}\mathbf{M}^T)}} \int_{\mathbf{M}\Theta} u(\mathbf{v}) d\mathbf{v} \quad (2.3)$$

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

Though quite natural, proving Lemma 2.3 rigorously requires nontrivial arguments from geometric measure theory; see Section 9.1.

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 [23] 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.

(P.1) The noise vector $\xi = (\xi_1, \dots, \xi_n)^\top$ has zero mean.

(P.2) For any $\gamma > 0$ small enough, there exist a probability space and two random variables ξ' and ζ defined on this probability space satisfying the three following points:

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

$$\limsup_{\gamma \rightarrow 0} \sup_{\substack{(\mathbf{t}, \mathbf{a}) \in \mathbb{R}^p \times \mathbb{R}^p \\ (\|\mathbf{t}\|_2, \mathbf{a}) \in [0, t_0] \times \text{supp}(\xi')}} \frac{\ln \mathbb{E}\{\exp(\mathbf{t}^\top \zeta) | \xi' = \mathbf{a}\}}{\|\mathbf{t}\|_2^2 \gamma v(\mathbf{a})} \leq 1.$$

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

Assumption (P.2) is based on [23, Assumption N], and can be shown to be fulfilled for a large class of noise.

Proposition 3.1. Assume that ξ has zero mean. Assumption (P.2) is fulfilled with $t_0 = \infty$ when

- ξ is a Gaussian random vector with covariance matrix Σ , with $t_0 = \infty$ and $v(\mathbf{a}) \equiv \|\Sigma\|$;
- ξ is a Laplace random vector with covariance Σ , with $t_0 < \sqrt{2/\|\Sigma\|}$ and $v(\mathbf{a}) \equiv \|\Sigma\|/(1 - t_0^2 \|\Sigma\|^2/2)$;
- ξ is a bounded symmetric random vector, i.e., $\Pr\{|\xi_i| \leq B_i\}$ for some $\mathbf{B} \in \mathbb{R}^n$, with $t_0 = \infty$ and $v(\mathbf{a}) = \|\mathbf{a}\|_2 \leq \|\mathbf{B}\|_2$.

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

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

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 (P.1) and (P.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

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

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

The proof of Theorem 3.1 is a mild adaptation of the original one in [23, Section 2], where we used directly Assumption (P.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 [19] proposed three types of oracle inequalities which are similar to (3.2) under different assumptions. The first type (see [19, Theorem 1]) holds under a restrictive condition on the noise. The second (see [19, Theorem 2]) involves conditions depending on the noise and also on the dictionary. The last (see [19, 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 (P.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 [20]:

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

where \widehat{f}_{λ} , $\lambda \in \Lambda$ are affine estimators satisfying some conditions imposed in [20, Theorem 1] which yield the definition of \widehat{r}_{λ} , $\lambda \in \Lambda$. 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 [22] 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. EWA estimator

4.1. Group-analysis 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 $q \geq p$. We partition the index set $\{1, \dots, q\}$ into L non-overlapping groups/blocks of indices $\{\mathcal{G}_{\ell}\}_{1 \leq \ell \leq L}$ such that

$$\mathcal{G} = \bigcup_{\ell=1}^L \mathcal{G}_{\ell} = \{1, \dots, q\} \quad \text{and} \quad \mathcal{G}_{\ell} \cap \mathcal{G}_k = \emptyset, \quad \forall \ell \neq k.$$

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

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

which is a norm for $s \geq 1$, with $\|\mathbf{v}\|_{\infty, \mathcal{G}} = \max_{\ell \in \{1, \dots, L\}} \|\mathbf{v}_{\mathcal{G}_{\ell}}\|_2$. It is a quasi-norm for $s \in]0, 1[$. $\|\mathbf{v}\|_{0, \mathcal{G}}$ 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.

(H.1) Let $\mathbf{D} \in \mathbb{R}^{p \times q}$, there exists $\widetilde{\boldsymbol{\theta}} \in \mathbb{R}^p$ such that $f_{\widetilde{\boldsymbol{\theta}}} = f$ and $\|\mathbf{D}^{\top} \widetilde{\boldsymbol{\theta}}\|_{0, \mathcal{G}} \ll n$.

In plain words, Assumption (H.1) says that the number of active groups of $\mathbf{D}^{\top} \widetilde{\boldsymbol{\theta}}$ 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 $(\mathbf{D}^{\top} \widetilde{\boldsymbol{\theta}})_{\mathcal{G}_{\ell}}$ are nearly zero. We also impose the following assumption on \mathbf{D} .

(H.2) \mathbf{D} is a frame with frame bounds $\nu \geq \kappa > 0$.

Let us now introduce some applications in literature in which our sparsity context is mentioned.

Example 4.1 (2-D piecewise constant image). Let $\boldsymbol{\theta}_0 \in \mathbb{R}^{\sqrt{p} \times \sqrt{p}}$ be a 2-D piecewise constant image. In this framework, a popular analysis operator is the isotropic total variation called D_{TV} (see [53]). Namely, let $D_c : \mathbb{R}^{\sqrt{p} \times \sqrt{p}} \rightarrow \mathbb{R}^{\sqrt{p} \times \sqrt{p}}$ and $D_r : \mathbb{R}^{\sqrt{p} \times \sqrt{p}} \rightarrow \mathbb{R}^{\sqrt{p} \times \sqrt{p}}$ the finite difference operators along, respectively, the columns and rows of an image, with Neumann boundary conditions. We define D_{TV} as

$$D_{\text{TV}} : \boldsymbol{\theta} \in \mathbb{R}^{\sqrt{p} \times \sqrt{p}} \mapsto \text{vec} \left((\text{vec}(D_r(\boldsymbol{\theta})), \text{vec}(D_c(\boldsymbol{\theta})))^{\top} \right)^{\top} \in \mathbb{R}^{2p}.$$

By vectorizing $\tilde{\theta}$ in Assumption (H.1), \mathbf{D}^\top can be considered as the matrix version of the linear operator D_{TV} , called \mathbf{D}_{TV} . Here, $\mathbf{D}_{\text{TV}} = [\mathbf{D}_r^\top \ \mathbf{D}_c^\top]^\top$ where $\mathbf{D}_r \in \mathbb{R}^{p \times p}$ (resp. $\mathbf{D}_c \in \mathbb{R}^{p \times p}$) is the matrix counterpart of D_r (resp. D_c). With Neumann boundary conditions, \mathbf{D}_r and \mathbf{D}_c are bijective implying injectivity of \mathbf{D}_{TV} . Thus, \mathbf{D} is a frame in view of Courant-Fisher theorem.

The isotropic total variation prior on θ promotes sparsity of $\text{vec}(\| [D_{\text{TV}}(\theta_0)]_{i,j} \|_2)_{1 \leq i,j \leq \sqrt{p}}$. 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 $D_{\text{TV}}(\theta_0)$ with groups of size 2 along the third dimension.

Example 4.2 (Signal with overlapping groups). Consider $\theta_0 \in \mathbb{R}^p$ generated from L groups which overlap. The analysis operator acts as a group extractor (see [13, 47]). In this framework, $\mathbf{D}\mathbf{D}^\top = \text{diag}(\mathbf{B}_1^\top \mathbf{B}_1, \dots, \mathbf{B}_L^\top \mathbf{B}_L)$ where $\mathbf{B}_\ell \in \mathbb{R}^{q_\ell \times p}$, for $\ell \in \{1, \dots, L\}$, is a countable collection of localization operators, and then $q = \sum_{\ell=1}^L q_\ell$. Since the localization operators are injective, $\mathbf{D}\mathbf{D}^\top$ is bijective. Thus \mathbf{D} is a frame in view of Courant-Fisher theorem.

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.

4.2. Choice of dictionary

Let $\mathbf{X} \in \mathbb{R}^{n \times p}$ where $X_{i,j} = f_j(x_i)$, where we recall that $f_j \in \mathcal{H}$ is a known function (atom) in the deterministic dictionary \mathcal{H} .

(H.3) \mathbf{X} is normalized such that all the diagonal entries of $\mathbf{X}^\top \mathbf{X}/n$ are 1.

Now, let us introduce our dictionary of aggregation:

$$\mathcal{F}_\Theta = \{f_\theta = \mathcal{L}(\sum_{j=1}^p \theta_j f_j) : \theta \in \Theta = \{\theta \in \mathbb{R}^p : \|\mathbf{D}^\top \theta\|_{a,\mathcal{G}}^a \leq R\}\}, \quad (4.1)$$

where $a \in]0, 1]$, $R \in]0, +\infty]$ and $\mathcal{L} : \mathbb{R} \rightarrow \mathbb{R}$ is twice continuously differentiable and known function that depends on the regression problem: e.g., $\mathcal{L}(x) = e^x$ for the exponential regression, $\mathcal{L}(x) = e^x/(e^x + 1)$ for the logistic regression and $\mathcal{L}(x) = x$ for the linear regression. This dictionary of aggregation is similar to the one proposed in [21–23]. However, the set of indices is modified to adapt to 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 (4.1) is established in the following result.

Proposition 4.1. Let $\mathcal{F}_\Theta = \{f_\theta : \theta \in \Theta\}$ defined in (4.1) with some $R > 0$, $a \in]0, 1]$ and $\mathcal{L} : \mathbb{R} \rightarrow \mathbb{R}$ twice continuously differentiable. Let Assumption (H.2) holds for some $\kappa > 0$. We get that

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

where $\mathcal{B} = \{x \in \mathbb{R}^n : \|x\|_2 \leq \|\mathbf{X}\| R^{1/a} \kappa^{-1/2}\}$ and $\mathcal{L} : x \in \mathbb{R}^n \rightarrow (\mathcal{L}(x_1), \dots, \mathcal{L}(x_n))$.

From Proposition 4.1, one can choose $H = 2 \max_{x \in \mathcal{B}} \|\mathcal{L}(x)\|_2$.

Remark 4.1. By choosing $H = 2 \max_{x \in \mathcal{B}} \|\mathcal{L}(x)\|_2$, H depends on \mathbf{X} and then on n under Assumption (H.3). So $\beta \geq \max(4 \|v\|_\infty, 2H/t_0)$ also depends on n . In this case, ξ must satisfy Assumption (P.2) with $t_0 = \infty$. In view of Proposition 3.1, we can consider ξ as a Gaussian or a bounded symmetric noise.

4.3. Choice of prior

4.3.1. Main assumptions

Recall that the goal is to find a prior leading an oracle inequality with a small remainder term while promoting group-analysis sparsity. In order to promote sparsity, it is well-known that the prior is usually expected to be symmetric and sharply peaked around its mode (the origin) while having tails heavier than merely exponential [44]. This is the rationale behind our general prior which takes the form

$$\pi(d\theta) = \frac{1}{C_{\alpha,g,R}} \prod_{\ell=1}^L \exp(-\alpha^\ell \|\mathbf{D}^\top \theta\|_{\mathcal{G}_\ell}^a) g(\|\mathbf{D}^\top \theta\|_{\mathcal{G}_\ell}) I_\Theta(\theta) d\theta, \quad (4.2)$$

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

- (G.1) Boundedness: $g : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ is a bounded function such that $g \not\equiv 0$, and $\theta \mapsto g(\|[\mathbf{D}^\top \theta]_{\mathcal{G}_\ell}\|_2)$ is measurable on \mathbb{R}^p , $\forall \ell \in \{1, \dots, L\}$.
- (G.2) Integrability: $\int_{\mathbb{R}^p} \prod_{\ell=1}^L g(\|[\mathbf{D}^\top \mathbf{u}]_{\mathcal{G}_\ell}\|_2) d\mathbf{u} < \infty$.
- (G.3) Moment condition: $\int_{\mathbb{R}^p} \|[\mathbf{D}^\top \mathbf{u}]_{\mathcal{G}_\ell}\|_2^2 \prod_{k=1}^L g(\|[\mathbf{D}^\top \mathbf{u}]_{\mathcal{G}_k}\|_2) d\mathbf{u} < \infty$, $\forall \ell \in \{1, \dots, L\}$.
- (G.4) Growth condition: there exist $\lambda \geq 0$ and $h : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ such that $g(\|\mathbf{t} - \mathbf{t}^*\|_2)/g(\|\mathbf{t}\|_2) \leq h(\|\mathbf{t}^*\|_2)^\lambda$, $\forall (\mathbf{t}, \mathbf{t}^*) \in \mathbb{R}^G \times \mathbb{R}^G$.

The exponential part can be viewed as a generalized Gaussian on the group analysis coefficients with shape parameter a and scale parameter α^{-1} . The choice of $a \in]0, 1]$ favors sparsity. The choice of α will be made clear when discussing the remainder terms in our group-analysis SOI (see for instance Remark 5.1). The role of the function g is (at least) twofold. The first is precisely to capture heavier tails than exponential. Second, it allows more flexibility to adjust to the group sparsity scenario and optimize the performance of EWA (see remarks hereafter for further discussion and examples).

4.3.2. Discussion of the assumptions

Assumption (G.4) controls the growth of the function g . The growth function h will impact the remainder term in the main group-analysis SOI stated in Theorem 5.1, and more precisely the term $\Omega_{\mu_n, n, L, \lambda}^{\mathbf{D}}(\theta)$ therein.

Assumptions (G.1)-(G.3) play a prominent role in controlling the key constant $K_{a,g}^{\mathbf{D}} > 0$ that is involved in the construction of our general group-analysis SOI in Theorem 5.1. The following remark formalizes the existence of this constant and relates it to the integrability and moment conditions on g .

Remark 4.2. Let $G \geq 1$ and $\mathbf{D} \in \mathbb{R}^{p \times q}$ satisfying Assumption (H.2). For any function g satisfying Assumptions (G.1)-(G.3), and any $a \in]0, 1]$, there exists $K_{a,g}^{\mathbf{D}} \in]0, \infty[$ such that, $\forall \ell \in \{1, \dots, L\}$,

$$\frac{\int_{\mathbb{R}^p} \|[\mathbf{D}^\top \mathbf{u}]_{\mathcal{G}_\ell}\|_2^{2a} \prod_{k=1}^L g(\|[\mathbf{D}^\top \mathbf{u}]_{\mathcal{G}_k}\|_2) d\mathbf{u}}{\int_{\mathbb{R}^p} \prod_{k=1}^L g(\|[\mathbf{D}^\top \mathbf{v}]_{\mathcal{G}_k}\|_2) d\mathbf{v}} \leq K_{a,g}^{\mathbf{D}}. \quad (4.3)$$

Proof. From Assumption (G.3) and the fact that $g \not\equiv 0$, one can show that (4.3) holds for $a = 1$. Moreover, since $g \not\equiv 0$ and g satisfies Assumption (G.2), it holds that

$$\mathbf{u} \mapsto \frac{\prod_{\ell=1}^L g(\|[\mathbf{D}^\top \mathbf{u}]_{\mathcal{G}_\ell}\|_2) d\mathbf{u}}{\int_{\mathbb{R}^p} \prod_{k=1}^L g(\|[\mathbf{D}^\top \mathbf{v}]_{\mathcal{G}_k}\|_2) d\mathbf{v}}$$

is a probability measure. Therefore, (4.3) holds for any a in $]0, 1]$ by Hölder's inequality. \square

It is then legitimate to wonder whether there exists a simpler condition on g which implies Assumptions (G.2) and (G.3). The answer is affirmative as we now state.

Lemma 4.1. Let $G \geq 1$ and $\mathbf{D} \in \mathbb{R}^{p \times q}$ satisfying Assumption (H.2). Suppose that g satisfies Assumption (G.1) and

$$\int_0^\infty z^{G+1} g(z) dz < \infty. \quad (4.4)$$

Then Assumptions (G.2) and (G.3) are in force.

When the operator \mathbf{D} is invertible, $K_{a,g}^{\mathbf{D}}$ takes an even simpler and explicit form, and moreover, (4.4) is necessary for g to obey Assumption (G.3). This is the subject of Remark 4.3 and Remark 4.4.

Remark 4.3. Let $G \geq 1$ and $\mathbf{D} \in \mathbb{R}^{p \times p}$ be invertible. For any function g satisfying Assumptions (G.1)-(G.3), and for any $a \in]0, 1]$, one can choose $K_{a,g}^{\mathbf{D}}$ in (4.3) as

$$K_{a,g}^{\mathbf{D}} = \frac{\int_0^\infty x^{G-1+2a} g(x) dx}{\int_0^\infty z^{G-1} g(z) dz}. \quad (4.5)$$

Proof. The proof follows by combining Lemmas 2.3 and 2.2, i.e.,

$$\begin{aligned} \frac{\int_{\mathbb{R}^p} \left\| [\mathbf{D}^\top \mathbf{u}]_{\mathcal{G}_\ell} \right\|_2^{2a} \prod_{k=1}^L g\left(\left\| [\mathbf{D}^\top \mathbf{u}]_{\mathcal{G}_k} \right\|_2\right) d\mathbf{u}}{\int_{\mathbb{R}^p} \prod_{k=1}^L g\left(\left\| [\mathbf{D}^\top \mathbf{v}]_{\mathcal{G}_k} \right\|_2\right) d\mathbf{v}} &= \frac{\int_{\mathbb{R}^G} \|\mathbf{u}\|_2^{2a} g(\|\mathbf{u}\|_2) d\mathbf{u} \left(\int_{\mathbb{R}^G} g(\|\mathbf{v}\|_2) d\mathbf{v} \right)^{L-1}}{\left(\int_{\mathbb{R}^G} g(\|\mathbf{w}\|_2) d\mathbf{w} \right)^L} \\ &= \frac{\int_{\mathbb{R}^G} \|\mathbf{u}\|_2^{2a} g(\|\mathbf{u}\|_2) d\mathbf{u}}{\int_{\mathbb{R}^G} g(\|\mathbf{v}\|_2) d\mathbf{v}} \\ &= \frac{\int_0^\infty x^{G-1+2a} g(x) dx}{\int_0^\infty z^{G-1} g(z) dz}. \end{aligned}$$

□

Remark 4.4. When \mathbf{D} is invertible, if g does not satisfy (4.4) then g cannot fulfill Assumption (G.3). Consequently, Assumption (G.3) and condition (4.4) are equivalent in the invertible case.

Proof. By Lemmas 2.3 and 2.2, we get

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

□

4.3.3. Examples of g

Let us now discuss some choices of g .

Example 4.3. 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 [23] for the individual sparsity scenario, i.e., with $\mathbf{D} = \mathbf{I}_p$, $G = 1$ and $a = 1$.

Example 4.4. 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 Assumptions (G.2) and (G.3). Thanks to the parameters b and c , this choice of g offers more flexibility than the one in the previous example. This allows for example to optimize the performance of EWA by tuning these parameters for the particular dataset at hand.

5. 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 (4.1) and the prior (4.2) with a function g obeying Assumptions (G.1)-(G.4).

Theorem 5.1 (General group-analysis sparse oracle inequality). *Let $G \geq 1$, X satisfying Assumption (H.3), and D satisfying Assumption (H.2) with $\kappa > 0$. Let Assumptions (P.1) and (P.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 (4.1) and the prior (4.2) with g satisfying Assumptions (G.1)-(G.4). Let $K_{a,g}^D$, as defined in (4.3), and assume that $R > 3\sqrt{K_{a,g}^D}L$. Then the following group-analysis SOI holds,*

$$\mathbb{E} \left\{ \|\widehat{f}_n - f\|_n^2 \right\} \leq \inf_{\Theta_{\mu_n, L, R}^D} \left(\|f_\theta - f\|_n^2 + \Phi_{\mu_n, n, L}^D(\theta) + \Omega_{\mu_n, n, L, \lambda}^D(\theta) \right) + \Psi_{\mu_n, L, p}^D, \quad (5.1)$$

with

$$\begin{cases} \Theta_{\mu_n, L, R}^D = \{\theta \in \mathbb{R}^p : \|D^\top \theta\|_{a, \mathcal{G}}^a \leq R - 3\sqrt{K_{a,g}^D}L\} \\ \Phi_{\mu_n, n, L}^D(\theta) = \frac{\beta}{n} \left(1 + 3\sqrt{K_{a,g}^D}L\alpha^a + \alpha^a \|D^\top \theta\|_{a, \mathcal{G}}^a \right) \\ \Omega_{\mu_n, n, L, \lambda}^D(\theta) = \frac{\lambda\beta}{n} \sum_{\ell=1}^L \ln h(\| [D^\top \theta]_{\mathcal{G}_\ell} \|_2) \\ \Psi_{\mu_n, L, p}^D = 2\kappa^{-1} K_{1,g}^D \exp\left(3\sqrt{K_{a,g}^D}L\alpha^a\right) p C_{f, \mathcal{L}}, \end{cases}$$

and $C_{f, \mathcal{L}} = \|\mathcal{L}'\|_\infty^2 + \|\mathcal{L}''\|_\infty (\|\mathcal{L}\|_\infty + \|f\|_\infty)$.

Before proceeding, we pause to make a few important remarks.

Remark 5.1. *The group-analysis SOI (5.1) is sharp. It depends on several parameters as discussed below.*

- (i) *The parameter R appears in the dictionary. Namely, the EWA estimator \widehat{f}_n mimics the best aggregate f_θ for all possible weights belonging to $\{\theta \in \mathbb{R}^p : \|D^\top \theta\|_{a, \mathcal{G}}^a \leq R - 3\sqrt{K_{a,g}^D}L\}$. Then R must be sufficiently large to cover the “good” model f_θ in Assumption (H.1). Moreover, since $R > 3\sqrt{K_{a,g}^D}L$, $R \sim \sqrt{K_{a,g}^D}L$ is the smallest choice we can make to reduce the rate of $\Phi_{\mu_n, n, L}^D(\theta)$ as $\|D^\top \theta\|_{a, \mathcal{G}}^a \leq R$.*
- (ii) *The parameter α is used to cancel the effect of $L\sqrt{K_{a,g}^D}$ in the remainder terms. By choosing $\alpha \leq (3K\sqrt{K_{a,g}^D})^{-1/a}$, $\Phi_{\mu_n, n, L}^D(\theta) \leq \beta n^{-1}(1 + 3\sqrt{K_{a,g}^D}\alpha^a + \alpha^a R) \sim n^{-1}$ and $\Psi_{\mu_n, L, p}^D \leq \frac{2eC_{f, \mathcal{L}}}{\kappa} K_{1,g}^D p$.*
- (iii) *The parameter $K_{a,g}^D$ and the function h depend on the choice of g . They respectively control the rate of $\Psi_{\mu_n, L, p}^D$ and $\Omega_{\mu_n, n, L, \lambda}^D(\theta)$.*

In what follows, let us state the consequences of Theorem 5.1 with the choices of g in Example 4.3 and 4.4. Especially, we will discuss the rate of $\Omega_{\mu_n, n, L, \lambda}^D(\theta)$ and $\Psi_{\mu_n, L, p}^D$.

We first consider the prior (4.2) in Example 4.3, under the individual sparsity scenario ($D = I_p$, $G = 1$) and the choice $a = 1$ (i.e., $\Theta = \{\theta \in \mathbb{R}^p : \|\theta\|_1 \leq R\}$). This is the setting considered in [23]. We obtain the following SOI as a corollary of our main result.

Corollary 5.1. *Let X satisfying Assumption (H.3), $D = I_p$ and fix $G = 1$. Suppose that Assumptions (P.1) and (P.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 (4.1) and the prior (4.2) with g defined in Example 4.3 and $\alpha \leq 1/(3p\tau)$. Assume that $R > 3p\tau$. Choosing $\tau^2 \sim 1/(pn)$ and $R \sim p\tau$, SOI (5.1) holds with $\Theta_{\mu_n, L, R}^D = \{\theta \in \mathbb{R}^p : \|\theta\|_1 \leq R - 3p\tau\}$,*

$$\Omega_{\mu_n, n, L, \lambda}^D(\theta) \sim \frac{\|\theta\|_0 \ln(p)}{n}, \quad \text{and} \quad \Psi_{\mu_n, L, p}^D \sim \frac{1}{n}.$$

The order of $\Omega_{\mu_n, n, L, \lambda}^D(\theta)$ is the classical rate under the sparsity scenario. This scaling is similar as in [23] with the same prior. 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 5.2. 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 $\int_0^\infty x^{G+1}(\tau^2 + x^2)^{-2} dx$ is not definite. Consequently, condition (4.4) is not fulfilled with g defined in Example 4.3 when $G \geq 2$.

According to Remark 4.4, Remark 5.2 implies that Assumption (G.3) is not fulfilled for g in Example 4.3 when the group size $G \geq 2$ and \mathbf{D} invertible. Thus one cannot construct a group-analysis SOI from Theorem 5.1 to guarantee the quality of the corresponding estimator. Overcoming this limitation was yet another motivation behind the choice of g in Example 4.4, which turns out to work well under the group-analysis sparsity scenario. In a nutshell, an aggregate with g in Example 4.4 exhibits the group-analysis SOI defined in the following corollary with any $G \geq 1$, any $\mathbf{D} \in \mathbb{R}^{p \times q}$ satisfying Assumption (H.2) and any $a \in]0, 1]$.

Corollary 5.2. Let \mathbf{X} satisfying Assumption (H.3), $G \geq 1$ and \mathbf{D} satisfying Assumption (H.2) with $\kappa > 0$. Let Assumptions (P.1) and (P.2) be satisfied with some function v , (3.1) holds and $\beta \geq \max(4\|v\|_\infty, 2H/t_0)$. Take the dictionary (4.1) and the prior (4.2) with $a \in]0, 1]$, $\alpha \geq 0$ and g defined in Example 4.4. We get that g satisfies Assumptions (G.1)-(G.4). Then, let $K_{a,g}^{\mathbf{D}}$ as defined in (4.3), and assume that $R > 3\sqrt{K_{a,g}^{\mathbf{D}}L}$. Then the group-analysis SOI (5.1) holds, with $\lambda = c$ and $h(x) = 1 + (x/\tau)^b$.

To get an explicit control of the remainder term, it is instructive to have a closed-form of $K_{a,g}^{\mathbf{D}}$. This can be done for instance when \mathbf{D} is invertible, see (4.5). The obtained group-analysis SOI is stated as follows.

Corollary 5.3. Consider the same framework as Corollary 5.2 with \mathbf{D} invertible. For $a \in]0, 1]$, let $\tilde{K}_{a,g}^{\mathbf{D}} = \frac{\Gamma((2a+G)/b)\Gamma(c-(2a+G)/b)}{\Gamma(G/b)\Gamma(c-G/b)}$, and set $\alpha \leq 1/\left(3\tau^a \sqrt{\tilde{K}_{a,g}^{\mathbf{D}}L}\right)^{1/a}$. Choosing $\tau^2 \sim 1/(pn)$ and $R \sim L\tau^a$, the group-analysis SOI (5.1) holds with $\Theta_{\mu_n, L, R}^{\mathbf{D}} = \{\boldsymbol{\theta} \in \mathbb{R}^p : \|\mathbf{D}^\top \boldsymbol{\theta}\|_{a, \mathcal{G}}^a \leq R - 3\tau^a \sqrt{\tilde{K}_{a,g}^{\mathbf{D}}L}\}$,

$$\Omega_{\mu_n, n, L, \lambda}^{\mathbf{D}}(\boldsymbol{\theta}) \sim \frac{\|\mathbf{D}^\top \boldsymbol{\theta}\|_{0, \mathcal{G}} \ln(L)}{n}, \quad \text{and} \quad \Psi_{\mu_n, L, p}^{\mathbf{D}} \sim \frac{1}{n}.$$

By Assumption (H.1), $\|\mathbf{D}^\top \boldsymbol{\theta}\|_{0, \mathcal{G}}$ is small when $\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}$ (with R must be sufficiently large to cover $\tilde{\boldsymbol{\theta}}$). Thus, $\|\mathbf{D}^\top \tilde{\boldsymbol{\theta}}\|_{0, \mathcal{G}} \ln(L)$ is small compared to n . Under the sparsity scenario, the order of $\Omega_{\mu_n, n, L, \lambda}^{\mathbf{D}}(\boldsymbol{\theta})$ becomes $O(\|\boldsymbol{\theta}\|_0 \ln(p)/n)$ which is the same rate as the aggregate with g in Example 4.3.

6. Forward-Backward proximal LMC algorithm

The goal of this section is to implement our EWA estimator with the probability measure (1.3) 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 $\mathcal{L}(x) = x$, and thus²

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

where $\mathbf{X} \in \mathbb{R}^{n \times p}$ is the design matrix (see Section 4.2). Recall the EWA estimator

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

with the measure

$$\mu_n(\boldsymbol{\theta}) \propto \exp\left(-\frac{\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2}{\beta}\right) \pi(\boldsymbol{\theta}). \quad (6.1)$$

²Generalization to non-linear link functions is possible from a practical point of view. However, in this case, the convergence guarantees of the discrete Langevin-based MCMC sampling scheme would be much more intricate even for a simple exponential link function. One of the main difficulties lies in that the data loss $\boldsymbol{\theta} \mapsto \|\mathbf{y} - \mathcal{L}(\mathbf{X}\boldsymbol{\theta})\|_2^2$ is not necessarily convex which prevents us from showing that the Langevin diffusion is geometrically ergodic. This is an open question that we leave to a future work.

Computing $\widehat{\theta}_n$ 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_n by constructing an appropriate Markov chain whose stationary distribution is μ_n , 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 [23], we here use the Langevin diffusion process; see [51].

6.1. The Langevin diffusion

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

$$dL(t) = \frac{1}{2}\psi(L(t))dt + dW(t), \quad t > 0, \quad L(0) = I_0, \quad (6.2)$$

where $\psi = \nabla \ln(\mu)$, μ is everywhere non-zero and a suitably smooth target density function on \mathbb{R}^p , W is a p -dimensional Brownian process and $I_0 \in \mathbb{R}^p$ is the initial value. Under mild assumptions, the SDE (6.2) has a unique strong solution and, moreover, $L(t)$ has a stationary distribution with density precisely μ [51, Theorem 2.1]. $L(t)$ is therefore interesting for sampling from μ . In particular, this opens the door to approximating integrals $\int_{\mathbb{R}^p} f(\theta)\mu(\theta)d\theta$ by the average value of a Langevin diffusion, i.e., $\frac{1}{T} \int_0^T f(L(t))dt$ for a large enough T . Under additional assumptions on μ and f in a proper functional class, the expected squared error of the approximation can be controlled [57].

Forward Euler discretization In practice, in simulating the diffusion sample path, we cannot follow exactly the dynamic defined by the SDE (6.2). Instead, we must discretize it. A popular discretization is given by the forward (Euler) scheme, which reads

$$L_{k+1} = L_k + \frac{\delta}{2}\psi(L_k) + \sqrt{\delta}Z_k, \quad t > 0, \quad L_0 = I_0,$$

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

$$\frac{\delta}{T} \sum_{k=0}^{\lfloor T/\delta \rfloor - 1} L_k.$$

It is then tempting to approximate $\widehat{\theta}_n$ by applying this discretization strategy to the Langevin diffusion with μ_n in (6.1) as the target density. However, quantitative consistency guarantees of this discretization require μ (hence ψ) to be sufficiently smooth, which limits their applicability in our context. To cope with this difficulty, a few works have recently proposed to replace $\ln(\mu)$ with a smoothed version (typically involving the Moreau-Yosida regularization/envelope, see Definition 2.2) [28, 29, 46]. In [29, 46] for instance, the authors proposed proximal-type algorithms to sample from possibly non-smooth log-concave densities μ using the forward Euler discretization and the Moreau-Yosida regularization. In [46]³, $-\ln(\mu)$ is replaced with its Moreau envelope, while in [29], it is assumed that $-\ln(\mu) = F + G$, 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 [12], though without any guarantees.

6.2. Forward-Backward proximal Langevin MC algorithms

Consider the prior π in (4.2) with g is given in Example 4.4 and $H = +\infty$. Then, μ_n is neither differentiable nor log-concave. To overcome these difficulties, we will exploit the structure of μ_n and some arguments from variational

³The author however applied it to problems where $-\ln(\mu) = F + G$. But the gradient of the Moreau envelope of a sum, which amounts to computing the proximity operator of $-\ln(\mu)$ does not have an easily implementable expression even if those of F and G do.

analysis [52]. For ease of notation, in the following, we denote with the same symbol the measure and its density with respect to the Lebesgue measure. Thus μ_n reads

$$\mu_n(\boldsymbol{\theta}) \propto \exp(-V(\boldsymbol{\theta})) \quad (6.3)$$

where $V \stackrel{\text{def}}{=} F_\beta + W_\lambda \circ \mathbf{D}^\top$, $F_\beta(\boldsymbol{\theta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2/\beta$ and $W_\lambda(\mathbf{u}) = \sum_{\ell=1}^L w_\lambda(\|\mathbf{u}_{\mathcal{G}_\ell}\|_2)$ with

$$w_\lambda : x \in [0, +\infty[\mapsto \alpha^a x^a + c \ln(\tau^b + x^b), \quad (6.4)$$

and w_λ is parameterized by $\lambda = (a, b, c, \alpha, \tau) \in]0, 1] \times]0, 1] \times]2 + G/b, +\infty[\times \mathbb{R}^+ \times \mathbb{R}^{+,*}$.

We start by collecting some important properties on the function w_λ and its proximal operator, as well as their implications. We denote $V_\gamma = F_\beta + (\gamma W_\lambda) \circ \mathbf{D}^\top$.

Lemma 6.1. *The function w_λ in (6.4) is bounded from below, increasing and continuously differentiable on $]0, +\infty[$. Fix $a = b = 1$. Then, for any $\gamma \in]0, \tau^2/c[$, the following holds,*

(i) $\text{prox}_{\gamma w_\lambda}$ is single-valued on $[0, +\infty[$ and is given by

$$\text{prox}_{\gamma w_\lambda}(x) = \begin{cases} 0 & \text{if } x \leq \gamma w_\lambda'(0^+), \\ t - \gamma w_\lambda'(\text{prox}_{\gamma w_\lambda}(x)) & \text{if } x > \gamma w_\lambda'(0^+). \end{cases}$$

(ii) For any $x \geq 0$, $0 \leq \text{prox}_{\gamma w_\lambda}(x) \leq x$.

(iii) $\text{prox}_{\gamma W_\lambda}(\mathbf{u}) = \left(\text{prox}_{\gamma w_\lambda}(\|\mathbf{u}_{\mathcal{G}_1}\|_2) \frac{\mathbf{u}_{\mathcal{G}_1}^\top}{\|\mathbf{u}_{\mathcal{G}_1}\|_2}, \dots, \text{prox}_{\gamma w_\lambda}(\|\mathbf{u}_{\mathcal{G}_L}\|_2) \frac{\mathbf{u}_{\mathcal{G}_L}^\top}{\|\mathbf{u}_{\mathcal{G}_L}\|_2} \right)^\top$.

(iv) $\text{prox}_{\gamma W_\lambda \circ \mathbf{D}^\top}$ is Lipschitz continuous.

(v) $\nabla(\gamma W_\lambda \circ \mathbf{D}^\top) = \gamma^{-1} \mathbf{D} \circ (\mathbf{I}_q - \text{prox}_{\gamma W_\lambda}) \circ \mathbf{D}^\top$. It is a uniformly bounded and Lipschitz continuous operator and ∇V_γ is a Lipschitz continuous mapping.

(vi) Assume that $\gamma \in]0, \min(\tau^2/c, \beta/(2\|\mathbf{X}\|^2))]$. Let $\mathbf{M}_\gamma \stackrel{\text{def}}{=} \mathbf{I}_p - (2\gamma/\beta)\mathbf{X}^\top \mathbf{X}$ which is a symmetric positive definite matrix. Then, $\nabla^{M_\gamma, \gamma} V = \gamma^{-1} \mathbf{M}_\gamma (\mathbf{I}_p - \text{prox}_{\gamma W_\lambda \circ \mathbf{D}^\top} (\mathbf{I}_p - \gamma \nabla F_\beta))$ which is a Lipschitz continuous mapping.

We now describe two Langevin MC (LMC) sampling algorithms, originally proposed in [40], that are based on forward-backward proximal splitting and establish their guarantees for the penalty w_λ . In the rest of the section, we will fix $a = b = 1$.

Semi-Forward-Backward LMC (Semi-FBLMC) Assume that $\gamma \in]0, \tau^2/c[$. Define the following SDE with the Moreau-Yosida regularized version of W_λ

$$d\mathbf{L}(t) = -\frac{1}{2} \nabla V_\gamma(\mathbf{L}(t)) dt + d\mathbf{W}(t), \quad t > 0. \quad (6.5)$$

Inserting Lemma 6.1(v) into (6.5), the Euler discretization of (6.5) reads

$$\mathbf{L}_{k+1} = \mathbf{L}_k - \frac{\delta}{2} \nabla F_\beta(\mathbf{L}_k) - \frac{\delta}{2\gamma} \mathbf{D}(\mathbf{D}^\top \mathbf{L}_k - \text{prox}_{\gamma W_\lambda}(\mathbf{D}^\top \mathbf{L}_k)) + \sqrt{\delta} \mathbf{Z}_k, \quad t > 0, \quad \mathbf{L}_0 = \mathbf{l}_0. \quad (6.6)$$

Forward-Backward LMC (FBLMC) Assume that $\gamma \in]0, \min(\tau^2/c, \beta/(2\|\mathbf{X}\|^2))]$. One can consider an alternative version of the SDE (6.5) with the Moreau-Yosida regularized version of V in the metric \mathbf{M}_γ (see Lemma 6.1(vi)), i.e.,

$$d\mathbf{L}(t) = -\frac{1}{2} \nabla(\mathbf{M}_\gamma, \gamma V \circ \mathbf{M}_\gamma^{-1/2})(\mathbf{L}(t)) dt + \mathbf{M}_\gamma^{1/2} d\mathbf{W}(t), \quad t > 0. \quad (6.7)$$

By the change of variable $U(t) = M_\gamma^{-1/2} L(t)$ we get the following SDE

$$dU(t) = -\frac{1}{2} M_\gamma^{-1} \nabla^{M_\gamma, \gamma} V(U(t)) dt + dW(t), \quad t > 0. \quad (6.8)$$

In view of Lemma 6.1(vi), the Euler discretization of (6.8) is given by

$$U_{k+1} = (1 - \frac{\delta}{2\gamma}) U_k + \frac{\delta}{2\gamma} \text{prox}_{\gamma W_\lambda \circ D^\top} (U_k - \gamma \nabla F_\beta(U_k)) + \sqrt{\delta} Z_k, \quad t > 0, \quad U_0 = I_0. \quad (6.9)$$

Remark 6.1. Observe that $\text{prox}_{\gamma W_\lambda \circ D^\top}$ in (6.9) is no separable in general. Owing to Assumption (H.2), one can show quite immediately that

$$\text{prox}_{\gamma W_\lambda \circ D^\top} = \tilde{D} \circ \text{prox}_{\gamma W_\lambda + \iota_{\text{Im}(D^\top)}}^{D^\top \tilde{D}} \circ D^\top,$$

where $\iota_{\text{Im}(D^\top)}(\mathbf{u}) = 0$ if $\mathbf{u} \in \text{Im}(D^\top)$ and $+\infty$ otherwise, and $D^\top \tilde{D}$ is indeed definite positive on $\text{Im}(D^\top)$. This means that, unless D is orthogonal, $W_\lambda \circ D^\top$ does not have an easy-to-compute expression, but rather necessitates to solve an optimization subproblem. Thus, from a computational perspective, for a general frame D , Semi-FBLMC is more efficient. For the case where D is invertible, one can operate a simple change of variable $\mathbf{u} = D^\top \theta$ and replace F_β with $F_\beta \circ D^{\top-1}$, and $W_\lambda \circ D^\top$ with W_λ . In this case, FBLMC can be applied efficiently.

6.3. Convergence guarantees

By virtue of Lemma 6.1(v)-(vi), standard results (see, e.g., [51, Theorem 2.1]) show that for any initial point $L(0)$ such that $\mathbb{E} \{\|L(0)\|_2^2\} < \infty$, (6.5) has a unique solution which is strongly Markovian, $\mathbb{E} \{\|L(t)\|_2^2\} < \infty$ for all $t > 0$, and L admits an (unique) invariant measure having the density μ_γ

$$\mu_\gamma(\theta) \propto \exp(-V_\gamma(\theta)).$$

The same also holds for (6.7) with the corresponding invariant measure. The following claim is a consequence of [40, Proposition 3.1] and Lemma 6.1. In the sequel, $\|\nu\|_{\text{TV}}$ stands for the total variation norm of a signed measure ν .

Proposition 6.1. Let μ_γ be the invariant measure of either (6.5) or (6.7). Then, $\|\mu_\gamma - \mu_n\|_{\text{TV}} \rightarrow 0$ as $\gamma \rightarrow 0$.

We consider the Semi-FBLMC discretization (6.6) of the Langevin diffusion (6.5). Let \tilde{L}_δ be the continuous-time extension of the scheme (6.6)

$$\tilde{L}_\delta(t) \stackrel{\text{def}}{=} L_0 - \frac{1}{2} \int_0^t \nabla V_\gamma(\tilde{L}(s)) ds + \int_0^t dW(s),$$

where $\tilde{L}(t) = L_k$ for $t \in [k\delta, (k+1)\delta[$. We write $\mathbf{P}_L^{t,\gamma}(I_0, \Omega) = \Pr\{L(t) \in \Omega | L(0) = I_0\}$ for all Borel sets Ω and initial condition I_0 . Similarly, we denote $\mathbf{P}_{\tilde{L}_\delta}^{t,\delta,\gamma}(I_0, \Omega) = \Pr\{\tilde{L}_\delta(t) \in \Omega | \tilde{L}_\delta(0) = I_0\}$. The superscripts stress the dependence on the parameters.

We also define the sample average

$$\bar{L}_{\delta,T,\gamma} = \frac{\delta}{T} \sum_{k=0}^{\lfloor T/\delta \rfloor - 1} L_k,$$

and the EWA estimate

$$\widehat{\theta}_\gamma = \int_{\mathbb{R}^p} \theta \mu_\gamma(\theta) d\theta.$$

Theorem 6.1. The following holds:

$$(i) \lim_{\gamma \rightarrow 0} \lim_{T \rightarrow +\infty} \lim_{\delta \rightarrow 0} \left\| \mathbf{P}_{\tilde{L}_\delta}^{T,\delta,\gamma}(I_0, \cdot) - \mu_n \right\|_{\text{TV}} = 0.$$

$$(ii) \text{ Suppose that } I_0 = 0, \text{ then } \lim_{T \rightarrow +\infty} \lim_{\delta \rightarrow 0} \mathbb{E} \left\{ \left\| \bar{L}_{\delta,T,\gamma} - \widehat{\theta}_\gamma \right\|_2 \right\} = 0.$$

To prove these claims, we first show that the Langevin diffusion in (6.5) is uniformly geometrically ergodic. Then we follow standard arguments, invoking Lemma 6.1, Proposition 6.1, the Girsanov formula and Pinsker inequality⁴. Similarly to Theorem 6.1, convergence guarantees of the FBLMC discretization scheme (6.9) can be established. We omit the details here for the sake of brevity.

⁴We have made no effort to sharpen the constants and the rates, and for instance their dependence on the dimension. This is beyond the scope of this paper.

7. 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

$$Y = X\theta_0 + \xi,$$

where ξ is the noise, X is the design matrix. $\theta_0 \in \mathbb{R}^p$ is the unknown regression vector of interest assumed to obey Assumption (H.1). Since $H = +\infty$, we then have to choose a distribution on the noise ξ such that β is independent of H . Such type of distributions is specified in [23, Section 2]. For our implementation, we assume ξ i.i.d. $\sim \mathcal{N}(0, \sigma^2)$. The noise level σ is chosen according to the simulated θ_0 .

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: $a = b = 1$, $c > (2 + G)/b$, $\tau \sim 1/(pn)$ and $\alpha \leq 1/\left(3\sqrt{K_{a,g}^D L}\right)^{1/a}$, where $K_{a,g}^D$ is given in Corollary 5.3. 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.

Following the philosophy of reproducible research, all the code implementing our EWA algorithm and reproducing the experiments of this paper are made publicly available for download at <https://github.com/luudytung/GroupAnalyseEWAToolbox>.

7.1. 1-D signal recovery under group sparsity

In this example, we set $D = I_p$, 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 θ_0 are equal to 1 and we denote $S = \|\theta_0\|_0$ the sparsity level of θ_0 . Two types of sparsity behavior are considered: individual sparsity where $G_{\theta_0} = 1$; group structured sparsity with $G_{\theta_0} = 4$. Besides, the positions of the non-zero/active entries (for $G_{\theta_0} = 1$) or groups (for $G_{\theta_0} = 4$) are chosen randomly uniformly on $\{1, \dots, p\}$.

The experiments are performed by fixing $p = 128$, and taking $S \in 2^{\{2, \dots, 7\}}$, $n \in 2^{\{3, \dots, 7\}}$, step-size $\delta = 4\sigma^2/(np)$ and integration time $T = 3500$. The parameters in the prior are chosen to minimize the remainder term in the oracle inequality (5.1). For each (S, n) , and each value of G_{θ_0} , $N_{\text{rep}} = 20$ instances of the problem suite (X, θ_0, 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

$$\pi_{S,n} = \frac{1}{N_{\text{rep}}} \sum_{j=1}^{N_{\text{rep}}} I\left(\left\|\widehat{\theta}_n^{(j,S,n)} - \theta_0^{(j,S,n)}\right\|_n \leq \epsilon\right), \quad (7.1)$$

where $\epsilon > 0$ (we choose $\epsilon = 0.4$) and $\widehat{\theta}_n^{(j,S,n)}$ (resp. $\theta_0^{(j,S,n)}$) corresponds to $\widehat{\theta}_n$ (resp. θ_0) in the j -th replication of (S, n) .

S/p and n/p are respectively normalized measures of sparsity and problem indeterminacy. We get a two-dimensional phase space $(S/p, n/p) \in [0, 1]^2$ describing the difficulty of a problem instance, i.e., problems are easier as one moves up (more measurements) and to the left (sparser θ_0). Phase diagrams plotting $\pi_{S,n}$ in (7.1) as a function $(S/p, n/p)$ were widely advocated by Donoho and co-authors for ℓ_1 minimization [24]. 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 that 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 $\widehat{\theta}_n^{(j,S,n)}$ when $(S, G_{\theta_0}) \in \{4, 8\} \times \{1, 4\}$ and $n/p = 1/2$. 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 $S/p = \|\theta_0\|_{0,\mathcal{G}} G_{\theta_0}/p$. As far as the expected phase transition curve is concerned, one has from Corollary 5.3 that it is expected to occur for

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

for some constant $C_\epsilon > 0$ depending on ϵ . That is, the phase transition curve is linear (p 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.

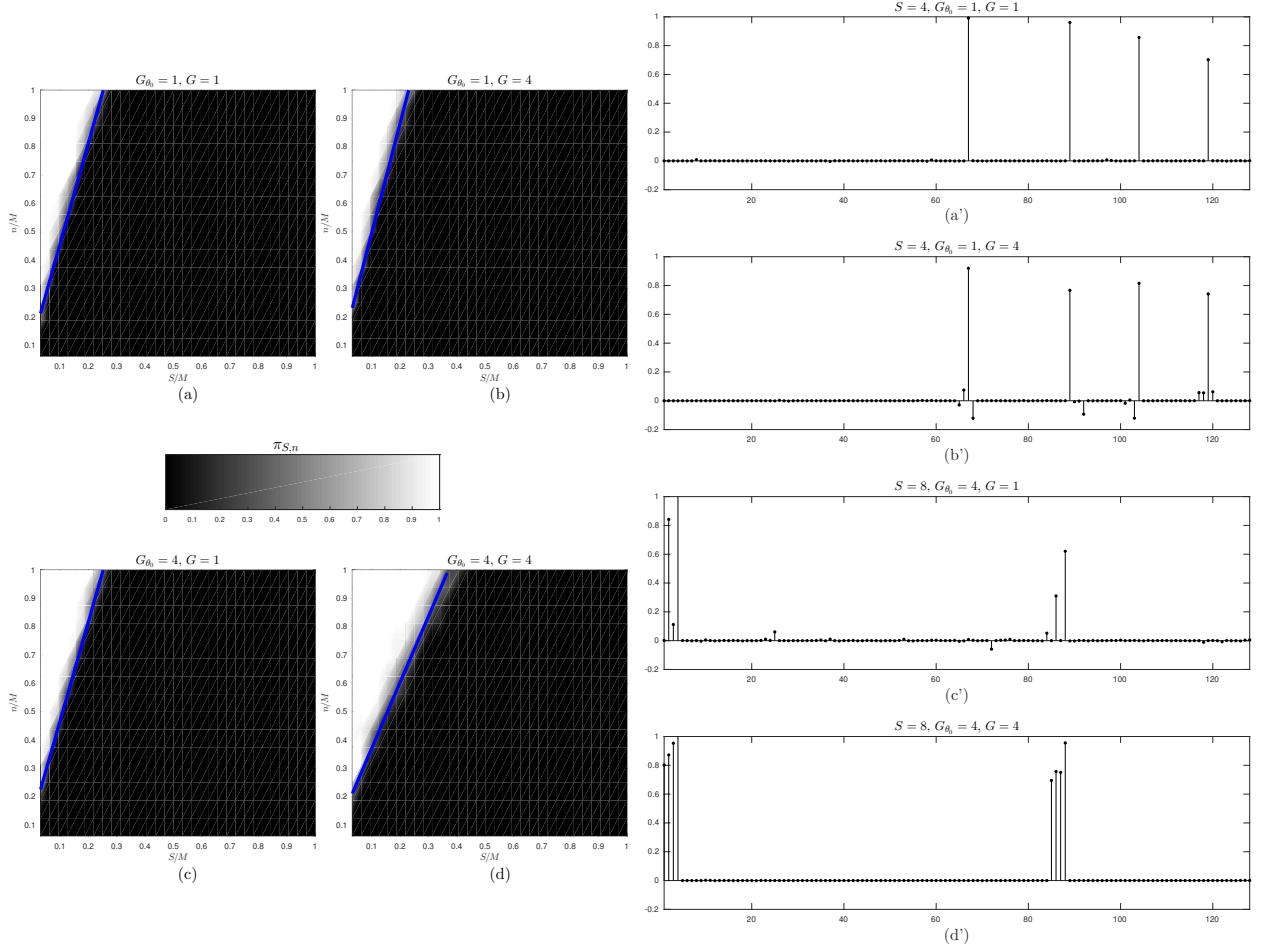


Figure 1: (a)-(d): Phase diagrams of EWA for $D = \mathbf{I}_p$, the color bar ranges from dark ($\pi_{S,n} = 0$) to bright ($\pi_{S,n} = 1$). The blue line is the fitted phase transition curve. (a')-(d'): Examples of vectors $\theta_n^{(j,S,n)}$ recovered by EWA with $n/p = 1/2$, two sparsity levels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

7.2. 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)). Thus $\text{vec}(\theta_0)$ is vector in \mathbb{R}^p with $p = 160^2$, and our goal is to recover θ_0 with values in $[0, 1]$ from

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

where $\xi \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$ and $X \in \mathbb{R}^{n \times p}$ 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 [53] which is

described in Example 4.1. It turns out that this can be cast in our analysis group-sparsity framework as a special case. In this experiment, we use Semi FBLMC Algorithm to compute the EWA estimator.

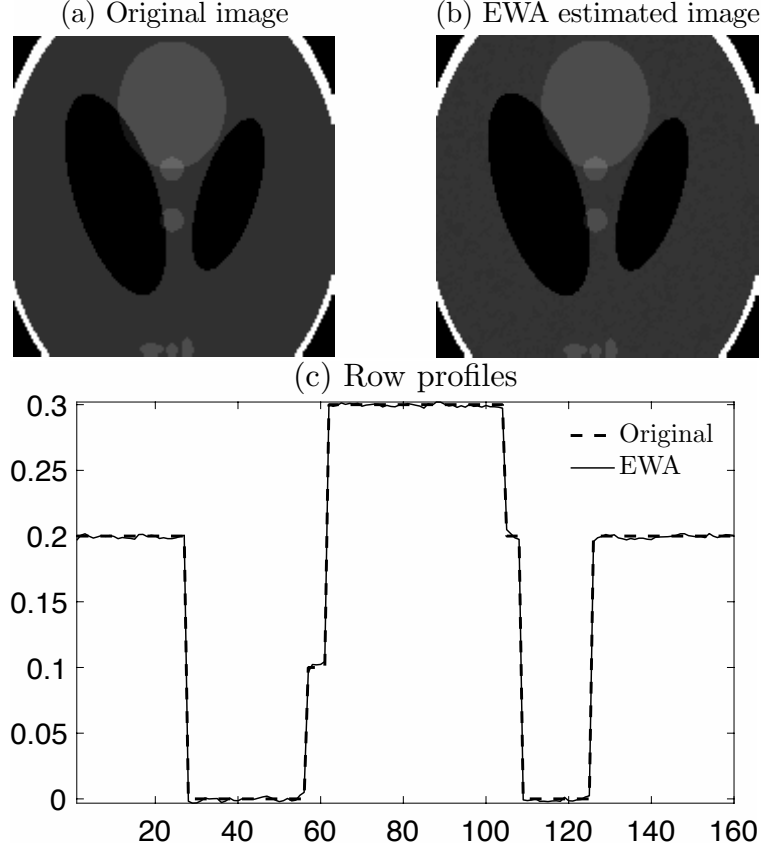


Figure 2: (a): Original close-up of Shepp-Logan phantom image. (b): Image recovered by EWA with $\delta = 2 \cdot 10^{-8}$ and $T = 10^4$. (c) Profiles of a row extracted from each image.

The results are depicted in Figure 2. $\sigma = 0.525$ in this experiment, the number of observations is $n = 9p/16 = 14400$, and we have $\|D_{TV}(\theta_0)\|_{0,\mathcal{G}} = 1376 \ll n$. A notable property of the EWA estimate is that it does not suffer from the stair-casing effect, unlike total variation minimization.

8. 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 convergence 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.

9. Proofs

9.1. Proofs of Section 2

Proof of Lemma 2.3 Consider the linear mapping $\mathbf{M} : \mathbf{x} \in \mathbb{R}^d \mapsto \mathbf{M}\mathbf{x} \in \mathbb{R}^r$, $r \geq d$. The Jacobian matrix of this mapping is obviously \mathbf{M} for any $\mathbf{x} \in \mathbb{R}^d$. Since \mathbf{M} is a frame, it is injective, hence so-called d -regular (see [41, Section 1.5]). In particular, $\det(\mathbf{M}\mathbf{M}^\top) > 0$. 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

$$\int_{\Theta} u(\mathbf{M}\mathbf{x})d\mathbf{x} = \frac{\int_{\mathbb{R}^r} \sum_{\mathbf{x} \in \Theta \cap \{\omega : \mathbf{M}\omega = \mathbf{v}\}} u(\mathbf{M}\mathbf{x})d\mathbf{v}}{\sqrt{\det(\mathbf{M}\mathbf{M}^\top)}} = \frac{\int_{\text{Im}(\mathbf{M})} \sum_{\mathbf{x} \in \Theta \cap \{\omega : \mathbf{M}\omega = \mathbf{v}\}} u(\mathbf{M}\mathbf{x})d\mathbf{v}}{\sqrt{\det(\mathbf{M}\mathbf{M}^\top)}}.$$

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

9.2. Proofs of Section 3

Proof of Proposition 3.1

- **Gaussian noise:** Let $\xi \sim \mathcal{N}(0, \Sigma)$. We set $\zeta \sim \mathcal{N}(0, (2\gamma + \gamma^2)\Sigma)$. Thus conditions (a) and (b) in Assumption (P.2) are verified. We check now condition (c). Let $\mathbf{t} \in \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_2 \in [-\infty, +\infty]\}$, $\mathbf{u} = (\mathbf{t}^\top \sqrt{2\gamma + \gamma^2} \Sigma^{1/2})^\top$ and $\epsilon \sim \mathcal{N}(0, \mathbf{I}_n)$, we get that

$$\mathbb{E}\{e^{\mathbf{t}^\top \zeta}\} = \mathbb{E}\{e^{\mathbf{u}^\top \epsilon}\} = \prod_{j=1}^n \mathbb{E}\{e^{u_j \epsilon_j}\} = e^{\frac{\|\mathbf{u}\|_2^2}{2}} \leq e^{\frac{1}{2} \|\sqrt{2\gamma + \gamma^2} \Sigma^{1/2} \mathbf{t}\|_2^2} \leq e^{(\gamma + \frac{\gamma^2}{2}) \|\Sigma\| \|\mathbf{t}\|_2^2}.$$

Thus, let $\mathbf{a} \in \mathbb{R}^n$ and $v(\mathbf{a}) \equiv \|\Sigma\|$, we then get

$$\frac{\ln \mathbb{E}\{e^{\mathbf{t}^\top \zeta} | \xi = \mathbf{a}\}}{\|\mathbf{t}\|_2^2 \gamma v(\mathbf{a})} \leq 1 + \frac{\gamma}{2} \xrightarrow{\gamma \rightarrow 0} 1 \leq 1.$$

- **Laplace noise:** Let $\xi \sim \mathcal{L}(0, \Sigma)$, i.e., its associated characteristic function is $\varphi_\xi(\mathbf{t}) = (1 + \mathbf{t}^\top \Sigma \mathbf{t}/2)^{-1}$, we choose ζ according to the distribution associated to the characteristic function

$$\varphi_\zeta(\mathbf{t}) = (1 + \gamma)^{-2} \left(1 + \frac{2\gamma + \gamma^2}{1 + (1 + \gamma)^2 \mathbf{t}^\top \Sigma \mathbf{t}/2} \right).$$

For any $\mathbf{t} \in \mathbb{R}^n$, we get that

$$\varphi_{\xi+\zeta}(\mathbf{t}) = \varphi_\xi(\mathbf{t})\varphi_\zeta(\mathbf{t}) = \frac{1}{1 + (1 + \gamma)^2 \mathbf{t}^\top \Sigma \mathbf{t}/2} = \varphi_{(1+\gamma)\xi}(\mathbf{t}). \quad (9.1)$$

Thus, $\zeta + \xi$ has the same distribution as $(1 + \gamma)\xi$. We also obtain $\mathbb{E}\{\zeta | \xi\} = \mathbb{E}\{\zeta\} = (-i)\nabla \varphi_\zeta(0) = 0$. It suffices to check condition (c) of Assumption (P.2). We know that

$$\mathbb{E}\{e^{\mathbf{t}^\top \zeta} | \xi\} = \mathbb{E}\{e^{\mathbf{t}^\top \zeta}\} = \varphi_\zeta(-i\mathbf{t}) = \frac{1}{(1+\gamma)^2} \left(1 + \frac{2\gamma + \gamma^2}{1 - (1 + \gamma)^2 \mathbf{t}^\top \Sigma \mathbf{t}/2} \right).$$

Using Taylor's formula, we have

$$\ln(\mathbb{E}\{e^{\mathbf{t}^\top \zeta} | \xi\}) = \frac{\gamma \mathbf{t}^\top \Sigma \mathbf{t}}{1 - \mathbf{t}^\top \Sigma \mathbf{t}/2} + O(\gamma^2).$$

Thus, let $\mathbf{t} \in \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_2 \leq t_0\}$, $\mathbf{a} \in \mathbb{R}^n$ and $v(\mathbf{a}) \equiv \frac{\|\Sigma\|}{1 - t_0^2 \|\Sigma\|^2/2}$, we get

$$\frac{\ln \mathbb{E}\{e^{\mathbf{t}^\top \zeta} | \xi = \mathbf{a}\}}{\|\mathbf{t}\|_2^2 \gamma v(\mathbf{a})} \xrightarrow{\gamma \rightarrow 0} \frac{\mathbf{t}^\top \Sigma \mathbf{t}/(1 - \mathbf{t}^\top \Sigma \mathbf{t}/2)}{\|\mathbf{t}\|_2^2 \|\Sigma\|/(1 - t_0^2 \|\Sigma\|^2/2)} \leq \frac{1 - t_0^2 \|\Sigma\|/2}{1 - \mathbf{t}^\top \Sigma \mathbf{t}/2} \leq \frac{1 - t_0^2 \|\Sigma\|/2}{1 - \|\mathbf{t}\|_2^2 \|\Sigma\|/2} \leq 1.$$

We get two last inequalities under the condition $1 - t_0^2 \|\Sigma\|/2 > 0$ equivalent $t_0 < \sqrt{2/\|\Sigma\|}$.

- **Bounded symmetric noise:** Let ξ are symmetric and $\Pr\{|\xi_i| \leq B_i\} = 1$ for some $\mathbf{B} \in \mathbb{R}^n$, we set $\zeta = (\zeta_1, \dots, \zeta_n)^\top$ such that $\zeta_i = (1 + \gamma)|\xi_i| \operatorname{sgn}(\operatorname{sgn}(\xi_i) - (1 + \gamma)U_i) - \xi_i$, $U_i \sim \mathcal{U}([-1, 1])$ for any $i \in \{1, \dots, n\}$. Using [23, Equation (22)], for any $\mathbf{t} \in \mathbb{R}^n$ and $\mathbf{a} \in \{\mathbf{x} \in \mathbb{R}^n : x_i \in [-B_i, B_i], \forall i \in \{1, \dots, n\}\}$, we get that

$$\mathbb{E}\{e^{\mathbf{t}^\top \xi} | \xi = \mathbf{a}\} = \prod_{j=1}^n \mathbb{E}\{e^{t_j \xi_j} | \xi_j = a_j\} = e^{-\mathbf{t}^\top \mathbf{a}} \left(e^{(1+\gamma)\mathbf{t}^\top \mathbf{a}} \frac{2+\gamma}{2+2\gamma} + e^{-(1+\gamma)\mathbf{t}^\top \mathbf{a}} \frac{\gamma}{2+2\gamma} \right). \quad (9.2)$$

From (9.2) and the symmetry of ξ , we obtain

$$\mathbb{E}\{e^{\mathbf{t}^\top (\zeta + \xi)}\} = \mathbb{E}\left\{\mathbb{E}\{e^{\mathbf{t}^\top (\zeta + \xi)} | \xi\}\right\} = e^{(1+\gamma)\mathbf{t}^\top \xi} \frac{2+\gamma}{2+2\gamma} + e^{-(1+\gamma)\mathbf{t}^\top \xi} \frac{\gamma}{2+2\gamma} = \mathbb{E}\{e^{(1+\gamma)\mathbf{t}^\top \xi}\}.$$

Thus, $\zeta + \xi$ has the same distribution as $(1 + \gamma)\xi$. Since $\mathbb{E}\{\zeta | \xi = \mathbf{a}\}$ equals to the gradient of $\mathbb{E}\{e^{\mathbf{t}^\top \xi} | \xi = \mathbf{a}\}$ at $\mathbf{t} = 0$, from (9.2) we have then $\mathbb{E}\{\zeta | \xi = \mathbf{a}\} = 0$, $\forall \mathbf{a} \in [-\mathbf{B}, \mathbf{B}]$. It suffices to check the condition (c) of Assumption (P.2). Owing to [19, Lemma 3] and [23, Equation (22)], we get that $\ln(\mathbb{E}\{e^{t_i \xi_i} | \xi_i = a_i\}) \leq (t_i a_i)^2 \gamma (1 + \gamma)$.

Thus, let $\mathbf{t} \in \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_2 \in [-\infty, \infty]\}$ and $v(\mathbf{a}) = \|\mathbf{a}\|_2^2$, we get that

$$\frac{\ln \mathbb{E}\{e^{\mathbf{t}^\top \xi} | \xi = \mathbf{a}\}}{\|\mathbf{t}\|_2^2 \gamma v(\mathbf{a})} = \frac{\sum_{i=1}^n \ln \mathbb{E}\{e^{t_i \xi_i} | \xi_i = a_i\}}{\|\mathbf{t}\|_2^2 \gamma \|\mathbf{a}\|_2^2} \leq \frac{\sum_{i=1}^n t_i^2 a_i^2 \gamma (1 + \gamma)}{\|\mathbf{t}\|_2^2 \gamma \|\mathbf{a}\|_2^2} \leq 1 + \gamma \xrightarrow{\gamma \rightarrow 0} 1 \leq 1.$$

9.3. Proofs of Section 4

Proof of Proposition 4.1 Let $\theta \in \Theta$ and $i \in \{1, \dots, n\}$. Setting $u_i^\theta = \sum_{j=1}^p \theta_j f_j(x_i)$ and $\mathbf{u}^\theta = (u_1^\theta, \dots, u_n^\theta)^\top$, and by virtue of (2.1), (2.2), (4.1) and the fact that $a \in]0, 1]$, we have

$$\|\mathbf{u}^\theta\|_2 = \|\mathbf{X}\theta\|_2 \leq \|\mathbf{X}\| \|\tilde{\mathbf{D}}\| \|\mathbf{D}^\top \theta\|_2 \leq \frac{\|\mathbf{X}\| \|\mathbf{D}^\top \theta\|_{a, \mathcal{G}}}{\sqrt{\kappa}} \leq \frac{\|\mathbf{X}\|_{R^{1/a}}}{\sqrt{\kappa}},$$

which in turn implies $\mathbf{u}^\theta \in \mathcal{B}$. Therefore, for any $(\theta, \theta') \in \Theta^2$, $\|\mathbf{f}_\theta - \mathbf{f}_{\theta'}\|_2 = \|\mathcal{L}(\mathbf{u}^\theta) - \mathcal{L}(\mathbf{u}^{\theta'})\|_2 \leq 2 \max_{\mathbf{x} \in \mathcal{B}} \|\mathcal{L}(\mathbf{x})\|_2$. \square

Proof of Lemma 4.1 Let us first check the integrability condition (G.2). By Lemmas 2.3 and 2.2, we obtain

$$\begin{aligned} \int_{\mathbb{R}^p} \prod_{\ell=1}^L g(\|\mathbf{D}^\top \mathbf{u}\|_{\mathcal{G}_\ell}) d\mathbf{u} &= \frac{\int_{\operatorname{Im}(\mathbf{D}^\top)} \prod_{\ell=1}^L g(\|\mathbf{v}_{\mathcal{G}_\ell}\|_2) d\mathbf{v}}{\sqrt{\det(\mathbf{D}\mathbf{D}^\top)}} \leq \frac{\int_{\mathbb{R}^q} \prod_{\ell=1}^L g(\|\mathbf{v}_{\mathcal{G}_\ell}\|_2) d\mathbf{v}}{\sqrt{\det(\mathbf{D}\mathbf{D}^\top)}} \\ &= \frac{\left(\int_{\mathbb{R}^G} g(\|\mathbf{u}\|_2) d\mathbf{u}\right)^L}{\sqrt{\det(\mathbf{D}\mathbf{D}^\top)}} \\ &= \frac{C_G^L \left(\int_0^\infty z^{G-1} g(z) dz\right)^L}{\sqrt{\det(\mathbf{D}\mathbf{D}^\top)}}. \end{aligned}$$

Since $G \geq 1$ and $g : \mathbb{R}^+ \rightarrow \mathbb{R}^+$, by (4.4), we get

$$\begin{aligned} \int_{\mathbb{R}^p} \prod_{\ell=1}^L g(\|\mathbf{D}^\top \mathbf{u}\|_{\mathcal{G}_\ell}) d\mathbf{u} &\leq \frac{C_G^L \left(\int_0^1 z^{G-1} g(z) dz + \int_1^\infty w^{G-1} g(w) dw\right)^L}{\sqrt{\det(\mathbf{D}\mathbf{D}^\top)}} \\ &\leq \frac{C_G^L \left(\sup_{z \in [0,1]} g(z) + \int_1^\infty w^{G+1} g(w) dw\right)^L}{\sqrt{\det(\mathbf{D}\mathbf{D}^\top)}} < \infty. \end{aligned} \quad (9.3)$$

Therefore, g satisfies Assumption (G.2). Now, we check the moment condition (G.3). Using similar arguments to the bound (9.3), we have

$$\begin{aligned}
\int_{\mathbb{R}^p} \left\| [\mathbf{D}^\top \mathbf{u}]_{\mathcal{G}_\ell} \right\|_2^2 \prod_{k=1}^L g(\left\| [\mathbf{D}^\top \mathbf{u}]_{\mathcal{G}_k} \right\|_2) d\mathbf{u} &\leq \frac{\int_{\mathbb{R}^q} \left\| \mathbf{v}_{\mathcal{G}_\ell} \right\|_2^2 \prod_{k=1}^L g(\left\| \mathbf{v}_{\mathcal{G}_k} \right\|_2) d\mathbf{v}}{\sqrt{\det(\mathbf{D}\mathbf{D}^\top)}} \\
&= \frac{\left(\int_{\mathbb{R}^G} \left\| \mathbf{u} \right\|_2^2 g(\left\| \mathbf{u} \right\|_2) d\mathbf{u} \right) \left(\int_{\mathbb{R}^G} g(\left\| \mathbf{v} \right\|_2) d\mathbf{v} \right)^{L-1}}{\sqrt{\det(\mathbf{D}\mathbf{D}^\top)}} \\
&= \frac{C_G^L \int_0^\infty z^{G+1} g(z) dz \left(\int_0^\infty w^{G-1} g(w) dw \right)^{L-1}}{\sqrt{\det(\mathbf{D}\mathbf{D}^\top)}} < \infty,
\end{aligned} \tag{9.4}$$

whence we conclude that g satisfies Assumption (G.3). \square

9.4. Proofs of Section 5

9.4.1. Proof of Theorem 5.1

Remind the prior $\pi(d\boldsymbol{\theta})$ from (4.2), where $\Theta = \{\boldsymbol{\theta} \in \mathbb{R}^p : \left\| \mathbf{D}^\top \boldsymbol{\theta} \right\|_{a,\mathcal{G}}^a \leq R\}$. Let $r_L = 3\sqrt{K_{a,g}^D L}$, $\Theta_{p_0^D} = \{\boldsymbol{\theta} \in \mathbb{R}^p : \left\| \mathbf{D}^\top \boldsymbol{\theta} - \mathbf{D}^\top \boldsymbol{\theta}^* \right\|_{a,\mathcal{G}}^a \leq r_L\}$ and

$$\boldsymbol{\theta}^* \in \{\boldsymbol{\theta} \in \mathbb{R}^p : \left\| \mathbf{D}^\top \boldsymbol{\theta} \right\|_{a,\mathcal{G}}^a \leq R - 3\sqrt{K_{a,g}^D L} = R - r_L\}. \tag{9.5}$$

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,

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

For any $i \in \{1, \dots, n\}$, with $\mathbf{X}_i = (f_1(x_i), \dots, f_p(x_i))^\top$, one can write $f_{\boldsymbol{\theta}}(x_i) = \mathcal{L}\left(\sum_{j=1}^p \theta_j f_j(x_i)\right) = \mathcal{L}(\mathbf{X}_i^\top \boldsymbol{\theta})$. Taylor-Lagrange formula then gives us

$$(f_{\boldsymbol{\theta}}(x_i) - f(x_i))^2 \leq (f_{\boldsymbol{\theta}^*}(x_i) - f(x_i))^2 + C_{f,\mathcal{L}} \left[\mathbf{X}_i^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \right]^2 + 2(f_{\boldsymbol{\theta}^*}(x_i) - f(x_i)) \mathcal{L}'(\mathbf{X}_i^\top \boldsymbol{\theta}^*) \mathbf{X}_i^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*), \tag{9.6}$$

where $C_{f,\mathcal{L}} = \left\| \mathcal{L}' \right\|_\infty^2 + \left\| \mathcal{L}'' \right\|_\infty (\left\| \mathcal{L} \right\|_\infty + \left\| f \right\|_\infty)$. By summing over i from 1 to n , normalizing by $1/n$, taking the integral in Θ w.r.t. p_0^D , inequality (9.6) becomes

$$\begin{aligned}
\int_{\Theta} \left\| f_{\boldsymbol{\theta}} - f \right\|_n^2 p_0^D(d\boldsymbol{\theta}) &\leq \left\| f_{\boldsymbol{\theta}^*} - f \right\|_n^2 + C_{f,\mathcal{L}} \int_{\mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \left[\mathbf{X}_i^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \right]^2 p_0^D(d\boldsymbol{\theta}) \\
&\quad + \frac{2}{n} \sum_{i=1}^n (f_{\boldsymbol{\theta}^*}(x_i) - f(x_i)) \mathcal{L}'(\mathbf{X}_i^\top \boldsymbol{\theta}^*) \mathbf{X}_i^\top \int_{\Theta} (\boldsymbol{\theta} - \boldsymbol{\theta}^*) p_0^D(d\boldsymbol{\theta}).
\end{aligned} \tag{9.7}$$

Note that, the right term of inequality (9.7) 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 last component vanishes. Indeed, let $\theta \in \Theta_{p_0^D}$, from (9.5) and the fact that $a \in]0, 1]$, we have

$$\begin{aligned}\|\mathbf{D}^\top \theta\|_{a, \mathcal{G}}^a &= \sum_{\ell=1}^L \|\mathbf{D}^\top \theta\|_{\mathcal{G}_\ell}^a \leq \sum_{\ell=1}^L \left(\|\mathbf{D}^\top \theta - \mathbf{D}^\top \theta^*\|_{\mathcal{G}_\ell}^a + \|\mathbf{D}^\top \theta^*\|_{\mathcal{G}_\ell}^a \right) \\ &\leq \|\mathbf{D}^\top \theta - \mathbf{D}^\top \theta^*\|_{a, \mathcal{G}}^a + \|\mathbf{D}^\top \theta^*\|_{a, \mathcal{G}}^a \\ &\leq r_L + \|\mathbf{D}^\top \theta^*\|_{a, \mathcal{G}}^a \leq R.\end{aligned}$$

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

$$\Theta_{p_0^D} \subseteq \Theta. \quad (9.8)$$

In what follows, we denote $\mathbb{B}_{a, \mathcal{G}}^a(x) = \{z \in \mathbb{R}^q : \|z\|_{a, \mathcal{G}}^a \leq x\}$, $\forall x > 0$ for brevity. By (9.8), property (2.1), Lemma 2.3 and symmetry of $\mathbb{B}_{a, \mathcal{G}}^a(r_L) \cap \text{Im}(\mathbf{D}^\top)$, we obtain

$$\begin{aligned}&\int_{\Theta} (\theta - \theta^*) p_0^D(d\theta) \\ &\propto \int_{\Theta \cap \Theta_{p_0^D}} (\theta - \theta^*) \prod_{\ell=1}^L \exp\left(-\alpha^a \|\mathbf{D}^\top \theta - \mathbf{D}^\top \theta^*\|_{\mathcal{G}_\ell}^a\right) g(\|\mathbf{D}^\top \theta - \mathbf{D}^\top \theta^*\|_{\mathcal{G}_\ell}) d\theta \\ &= \int_{\Theta_{p_0^D}} (\theta - \theta^*) \prod_{\ell=1}^L \exp\left(-\alpha^a \|\mathbf{D}^\top \theta - \mathbf{D}^\top \theta^*\|_{\mathcal{G}_\ell}^a\right) g(\|\mathbf{D}^\top \theta - \mathbf{D}^\top \theta^*\|_{\mathcal{G}_\ell}) d\theta \\ &= \frac{\tilde{\mathbf{D}}}{\sqrt{\det(\mathbf{D}\mathbf{D}^\top)}} \int_{\mathbb{B}_{a, \mathcal{G}}^a(r_L) \cap \text{Im}(\mathbf{D}^\top)} z \prod_{\ell=1}^L \exp\left(-\alpha^a \|z_{\mathcal{G}_\ell}\|_2^a\right) g(\|z_{\mathcal{G}_\ell}\|_2) dz = 0,\end{aligned} \quad (9.9)$$

which is the desired claim.

We now bound the second term in the right hand side of (9.7). Define

$$p_0(du) = \frac{1}{C_L \sqrt{\det(\mathbf{D}\mathbf{D}^\top)}} \prod_{\ell=1}^L \exp\left(-\alpha^a \|\mathbf{u}_{\mathcal{G}_\ell}\|_2^a\right) g(\|\mathbf{u}_{\mathcal{G}_\ell}\|_2) I_{\text{Im}(\mathbf{D}^\top) \cap \mathbb{B}_{a, \mathcal{G}}^a(r_L)}(\mathbf{u}) du. \quad (9.10)$$

One can see that p_0 coincides with the probability measure p_0^D on \mathbb{R}^p via a change of variables of type (2.3). So, p_0 is a probability measure on \mathbb{R}^q . For any $i, j \in \{1, \dots, L\}$, $i \neq j$, by a change of variables, we get $\int_{\mathbb{R}^q} \mathbf{u}_{\mathcal{G}_i} \mathbf{u}_{\mathcal{G}_j}^\top p_0(du) = - \int_{\mathbb{R}^q} \mathbf{u}_{\mathcal{G}_i} \mathbf{u}_{\mathcal{G}_j}^\top p_0(du)$, so

$$\int_{\mathbb{R}^q} \mathbf{u}_{\mathcal{G}_i} \mathbf{u}_{\mathcal{G}_j}^\top p_0(du) = 0. \quad (9.11)$$

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

$$\int_{\mathbb{R}^q} \mathbf{u}_{\mathcal{G}_j} \mathbf{u}_{\mathcal{G}_j}^\top p_0(du) = \int_{\mathbb{R}^q} \mathbf{u}_{\mathcal{G}_1} \mathbf{u}_{\mathcal{G}_1}^\top p_0(du). \quad (9.12)$$

We obtain

$$\begin{aligned}
& \int_{\mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \left[\mathbf{X}_i^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \right]^2 p_0^D(d\boldsymbol{\theta}) \\
& \stackrel{((2.1) \text{ and Lemma 2.3})}{=} \frac{1}{n} \int_{\mathbb{R}^q} \left[\mathbf{X} \tilde{\mathbf{D}} \mathbf{u} \right]^\top \mathbf{X} \tilde{\mathbf{D}} \mathbf{u} p_0(du) \\
& = \frac{1}{n} \int_{\mathbb{R}^q} \text{tr} \left(\mathbf{u}^\top \tilde{\mathbf{D}}^\top \mathbf{X}^\top \mathbf{X} \tilde{\mathbf{D}} \mathbf{u} \right) p_0(du) \\
& = \frac{1}{n} \text{tr} \left(\left(\mathbf{X} \tilde{\mathbf{D}} \right)^\top \mathbf{X} \tilde{\mathbf{D}} \int_{\mathbb{R}^q} \mathbf{u} \mathbf{u}^\top p_0(du) \right) \\
& \stackrel{((9.11) \text{ and } (9.12))}{=} \frac{1}{n} \sum_{\ell=1}^L \text{tr} \left(\left[\left(\mathbf{X} \tilde{\mathbf{D}} \right)^\top \mathbf{X} \tilde{\mathbf{D}} \right]_{\mathcal{G}_\ell} \int_{\mathbb{R}^q} \mathbf{u}_{\mathcal{G}_1} \mathbf{u}_{\mathcal{G}_1}^\top p_0(du) \right) \\
& \stackrel{(\text{Von Neumann's trace inequality})}{\leq} \frac{1}{n} \sum_{\ell=1}^L \sum_{j=1}^G \sigma_j \left(\left[\left(\mathbf{X} \tilde{\mathbf{D}} \right)^\top \mathbf{X} \tilde{\mathbf{D}} \right]_{\mathcal{G}_\ell} \right) \sigma_j \left(\int_{\mathbb{R}^q} \mathbf{u}_{\mathcal{G}_1} \mathbf{u}_{\mathcal{G}_1}^\top p_0(du) \right) \\
& \leq \frac{1}{n} \int_{\mathbb{R}^q} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 p_0(du) \sum_{\ell=1}^L \text{tr} \left(\left[\left(\mathbf{X} \tilde{\mathbf{D}} \right)^\top \mathbf{X} \tilde{\mathbf{D}} \right]_{\mathcal{G}_\ell} \right) \\
& = \frac{1}{n} \text{tr} \left(\left(\mathbf{X} \tilde{\mathbf{D}} \right)^\top \mathbf{X} \tilde{\mathbf{D}} \right) \int_{\mathbb{R}^q} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 p_0(du). \tag{9.13}
\end{aligned}$$

Moreover, by inequality (2.2), Assumption (H.3) and Von Neumann's trace inequality, we obtain

$$\frac{\text{tr} \left(\left(\mathbf{X} \tilde{\mathbf{D}} \right)^\top \mathbf{X} \tilde{\mathbf{D}} \right)}{n} \leq \sum_{j=1}^p \sigma_j \left(\frac{\mathbf{X}^\top \mathbf{X}}{n} \right) \sigma_j \left(\tilde{\mathbf{D}} \tilde{\mathbf{D}}^\top \right) \leq \sigma_1 \left(\tilde{\mathbf{D}} \tilde{\mathbf{D}}^\top \right) \sum_{j=1}^p \sigma_j \left(\frac{\mathbf{X}^\top \mathbf{X}}{n} \right) \leq \frac{p}{\kappa}. \tag{9.14}$$

Putting together (9.13) and (9.14), we get the bound

$$C_{f, \mathcal{L}} \int_{\mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \left[\mathbf{X}_i^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \right]^2 p_0^D(d\boldsymbol{\theta}) \leq C_{f, \mathcal{L}} \frac{p}{\kappa} \int_{\mathbb{R}^q} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 p_0(du). \tag{9.15}$$

Thanks to (9.9) and (9.15), inequality (9.7) becomes

$$\int_{\Theta} \|f_{\boldsymbol{\theta}} - f\|_n^2 p_0^D(d\boldsymbol{\theta}) \leq \|f_{\boldsymbol{\theta}^*} - f\|_n^2 + C_{f, \mathcal{L}} \frac{p}{\kappa} \int_{\mathbb{R}^q} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 p_0(du). \tag{9.16}$$

Now, inserting (9.16) into Theorem 3.1 (with $p = p_0^D$), we arrive at

$$\mathbb{E} \left\{ \|\widehat{f}_n - f\|_n^2 \right\} \leq \|f_{\boldsymbol{\theta}^*} - f\|_n^2 + C_{f, \mathcal{L}} \frac{p}{\kappa} \int_{\mathbb{R}^q} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 p_0(du) + \frac{\beta \text{KL}(p_0^D, \pi)}{n}. \tag{9.17}$$

To complete the proof, it remains to bound the last two terms in the right hand side of (9.17). This is the goal of the following lemma.

Lemma 9.1. *Consider the same framework as the one in Theorem 5.1, we have*

$$\int_{\mathbb{R}^q} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 p_0(du) \leq 2K_{1,g}^D e^{r_L \alpha^a}, \tag{9.18}$$

and

$$\text{KL}(p_0^D, \pi) \leq 1 + r_L \alpha^a + \lambda \sum_{\ell=1}^L \ln \left\{ h \left(\left\| \left[\mathbf{D}^\top \boldsymbol{\theta}^* \right]_{\mathcal{G}_\ell} \right\|_2 \right) \right\} + \alpha^a \left\| \mathbf{D}^\top \boldsymbol{\theta}^* \right\|_{a, \mathcal{G}}^a. \tag{9.19}$$

With $r_L = 3 \sqrt{K_{a,g}^D} L$, it follows from (9.17) and Lemma 9.1 that

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

According to (9.5), this completes the proof of Theorem 5.1. □

Proof of Lemma 9.1 To prove Lemma 9.1, we need an intermediate result.

Lemma 9.2. Let $s > L\sqrt{K_{a,g}^D}$. The following inequality holds

$$\frac{1}{T} \int_{\{u \in \mathbb{R}^p : \|\mathbf{D}^\top u\|_{a,g}^a > s\}} \prod_{\ell=1}^L g(\|\mathbf{D}^\top u\|_{\mathcal{G}_\ell}) du \leq \frac{L^2 K_{a,g}^D}{(s - L\sqrt{K_{a,g}^D})^2},$$

where $T = \int_{\mathbb{R}^p} \prod_{\ell=1}^L g(\|\mathbf{D}^\top u\|_{\mathcal{G}_\ell}) du$.

Proof of Lemma 9.2. Let \mathbf{U} be a random vector in \mathbb{R}^p with density $u \mapsto \frac{1}{T} \prod_{\ell=1}^L g(\|\mathbf{D}^\top u\|_{\mathcal{G}_\ell})$, where $T < \infty$ by Assumption (G.2). By Chebyshev inequality, we have

$$\begin{aligned} & \frac{1}{T} \int_{\{u \in \mathbb{R}^p : \|\mathbf{D}^\top u\|_{a,g}^a > s\}} \prod_{\ell=1}^L g(\|\mathbf{D}^\top u\|_{\mathcal{G}_\ell}) du \\ &= \Pr \left\{ \sum_{\ell=1}^L \|\mathbf{D}^\top \mathbf{U}\|_{\mathcal{G}_\ell}^a > s \right\} \\ &= \Pr \left\{ \sum_{\ell=1}^L \|\mathbf{D}^\top \mathbf{U}\|_{\mathcal{G}_\ell}^a - \mathbb{E} \left\{ \sum_{\ell=1}^L \|\mathbf{D}^\top \mathbf{U}\|_{\mathcal{G}_\ell}^a \right\} > s - \sum_{\ell=1}^L \mathbb{E} \left\{ \|\mathbf{D}^\top \mathbf{U}\|_{\mathcal{G}_\ell}^a \right\} \right\} \\ &\leq \mathbb{E} \left\{ \left[\sum_{\ell=1}^L \|\mathbf{D}^\top \mathbf{U}\|_{\mathcal{G}_\ell}^a - \mathbb{E} \left\{ \sum_{\ell=1}^L \|\mathbf{D}^\top \mathbf{U}\|_{\mathcal{G}_\ell}^a \right\} \right]^2 \right\} / (s - \sum_{\ell=1}^L \mathbb{E} \left\{ \|\mathbf{D}^\top \mathbf{U}\|_{\mathcal{G}_\ell}^a \right\})^2 \\ &= \text{var} \left(\sum_{\ell=1}^L \|\mathbf{D}^\top \mathbf{U}\|_{\mathcal{G}_\ell}^a \right) / (s - \sum_{\ell=1}^L \mathbb{E} \left\{ \|\mathbf{D}^\top \mathbf{U}\|_{\mathcal{G}_\ell}^a \right\})^2 \\ &\leq \mathbb{E} \left\{ \left(\sum_{\ell=1}^L \|\mathbf{D}^\top \mathbf{U}\|_{\mathcal{G}_\ell}^a \right)^2 \right\} / (s - \sum_{\ell=1}^L \mathbb{E} \left\{ \|\mathbf{D}^\top \mathbf{U}\|_{\mathcal{G}_\ell}^a \right\})^2. \end{aligned} \quad (9.20)$$

Next, by Cauchy-Schwartz inequality and Remark 4.2, we obtain

$$\mathbb{E} \left\{ \left(\sum_{\ell=1}^L \|\mathbf{D}^\top \mathbf{U}\|_{\mathcal{G}_\ell}^a \right)^2 \right\} \leq \mathbb{E} \left\{ L \sum_{\ell=1}^L \|\mathbf{D}^\top \mathbf{U}\|_{\mathcal{G}_\ell}^{2a} \right\} \leq L^2 K_{a,g}^D \quad (9.21)$$

and by Jensen inequality

$$s - \sum_{\ell=1}^L \mathbb{E} \left\{ \|\mathbf{D}^\top \mathbf{U}\|_{\mathcal{G}_\ell}^a \right\} \geq s - \sum_{\ell=1}^L \sqrt{\mathbb{E} \left\{ \|\mathbf{D}^\top \mathbf{U}\|_{\mathcal{G}_\ell}^{2a} \right\}} \geq s - L\sqrt{K_{a,g}^D} > 0. \quad (9.22)$$

Thus, combining (9.20), (9.21) and (9.22), we get

$$\frac{1}{T} \int_{\{u \in \mathbb{R}^p : \|\mathbf{D}^\top u\|_{a,g}^a > s\}} \prod_{\ell=1}^L g(\|\mathbf{D}^\top u\|_{\mathcal{G}_\ell}) du \leq \frac{L^2 K_{a,g}^D}{(s - L\sqrt{K_{a,g}^D})^2}.$$

□

We now turn to the proof of Lemma 9.1

Proof of Lemma 9.1. Let us begin by the proof of inequality (9.18). We have

$$\begin{aligned} \int_{\mathbb{R}^q} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 p_0(du) &= \frac{1}{C_L \sqrt{\det(\mathbf{D}\mathbf{D}^\top)}} \int_{\mathbb{B}_{a,\mathcal{G}}^a(r_L) \cap \text{Im}(\mathbf{D}^\top)} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 \prod_{\ell=1}^L e^{-\alpha^\ell \|\mathbf{u}_{\mathcal{G}_\ell}\|_2^a} g(\|\mathbf{u}_{\mathcal{G}_\ell}\|_2) du \\ &\leq \frac{1}{C_L \sqrt{\det(\mathbf{D}\mathbf{D}^\top)}} \int_{\mathbb{B}_{a,\mathcal{G}}^a(r_L) \cap \text{Im}(\mathbf{D}^\top)} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 \prod_{\ell=1}^L g(\|\mathbf{u}_{\mathcal{G}_\ell}\|_2) du. \end{aligned} \quad (9.23)$$

In the following, we show inequality (9.18) by bounding the right term of inequality (9.23). By Lemma 2.3 and Remark 4.2, we get

$$\begin{aligned} \frac{\int_{\mathbb{B}_{a,g}^a(r_L) \cap \text{Im}(\mathbf{D}^\top)} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 \prod_{\ell=1}^L g(\|\mathbf{u}_{\mathcal{G}_\ell}\|_2) d\mathbf{u}}{\sqrt{\det(\mathbf{D}\mathbf{D}^\top)} \int_{\mathbb{R}^p} \prod_{\ell=1}^L g(\|\mathbf{D}^\top \mathbf{u}\|_{\mathcal{G}_\ell}) d\mathbf{u}} &\leq \frac{\int_{\text{Im}(\mathbf{D}^\top)} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 \prod_{\ell=1}^L g(\|\mathbf{u}_{\mathcal{G}_\ell}\|_2) d\mathbf{u}}{\sqrt{\det(\mathbf{D}\mathbf{D}^\top)} \int_{\mathbb{R}^p} \prod_{\ell=1}^L g(\|\mathbf{D}^\top \mathbf{u}\|_{\mathcal{G}_\ell}) d\mathbf{u}} \\ &= \frac{\int_{\mathbb{R}^p} \|\mathbf{D}^\top \mathbf{u}\|_{\mathcal{G}_1}^2 \prod_{\ell=1}^L g(\|\mathbf{D}^\top \mathbf{u}\|_{\mathcal{G}_\ell}) d\mathbf{u}}{\int_{\mathbb{R}^p} \prod_{\ell=1}^L g(\|\mathbf{D}^\top \mathbf{u}\|_{\mathcal{G}_\ell}) d\mathbf{u}} \leq K_{1,g}^D. \end{aligned}$$

Then

$$\frac{1}{\sqrt{\det(\mathbf{D}\mathbf{D}^\top)}} \int_{\mathbb{B}_{a,g}^a(r_L) \cap \text{Im}(\mathbf{D}^\top)} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 \prod_{\ell=1}^L g(\|\mathbf{u}_{\mathcal{G}_\ell}\|_2) d\mathbf{u} \leq K_{1,g}^D T. \quad (9.24)$$

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

$$\begin{aligned} C_L^{-1} &= \left(\int_{\Theta_{p_0^D}} \prod_{\ell=1}^L e^{-\alpha^\ell} \|\mathbf{D}^\top \boldsymbol{\theta} - \mathbf{D}^\top \boldsymbol{\theta}^*\|_{\mathcal{G}_\ell}^2 g(\|\mathbf{D}^\top \boldsymbol{\theta} - \mathbf{D}^\top \boldsymbol{\theta}^*\|_{\mathcal{G}_\ell}) d\boldsymbol{\theta} \right)^{-1} \\ &= \left(\int_{\{u \in \mathbb{R}^p : \|\mathbf{D}^\top u\|_{a,g}^a \leq r_L\}} e^{-\alpha^\ell} \|\mathbf{D}^\top u\|_{a,g}^a \prod_{\ell=1}^L g(\|\mathbf{D}^\top u\|_{\mathcal{G}_\ell}) du \right)^{-1} \\ &\leq e^{r_L \alpha^a} \left(\int_{\{u \in \mathbb{R}^p : \|\mathbf{D}^\top u\|_{a,g}^a \leq r_L\}} \prod_{\ell=1}^L g(\|\mathbf{D}^\top u\|_{\mathcal{G}_\ell}) du \right)^{-1}. \end{aligned}$$

Since $r_L = 3\sqrt{K_{a,g}^D}L > \sqrt{K_{a,g}^D}L$, Lemma 9.2 gives us

$$\begin{aligned} C_L^{-1} &\leq e^{r_L \alpha^a} \left[T \left(1 - \frac{1}{T} \int_{\{u \in \mathbb{R}^p : \|\mathbf{D}^\top u\|_{a,g}^a > r_L\}} \prod_{\ell=1}^L g(\|\mathbf{D}^\top u\|_{\mathcal{G}_\ell}) du \right) \right]^{-1} \\ &\leq e^{r_L \alpha^a} T^{-1} \left(1 - \frac{L^2 K_{a,g}^D}{(r_L - L\sqrt{K_{a,g}^D})^2} \right)^{-1} \\ &= e^{r_L \alpha^a} T^{-1} \left(1 - \frac{1}{4} \right)^{-1} \leq 2e^{r_L \alpha^a} T^{-1}. \end{aligned} \quad (9.25)$$

Combining (9.24) and (9.25), (9.23) becomes $\int_{\mathbb{R}^q} \|\mathbf{u}_{\mathcal{G}_1}\|_2^2 p_0(d\mathbf{u}) \leq 2K_{1,g}^D e^{r_L \alpha^a}$. That concludes the proof of inequality (9.18) in Lemma 9.1.

Next, we prove inequality (9.19). Remind that $\text{supp}(\pi) = \Theta$, $\text{supp}(p_0^D) = \Theta_{p_0^D}$. By (9.8), 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 (9.19)

can be proved as follows. By Lemma 2.3, we have

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

where p_0 is a probability measure in \mathbb{R}^q defined in (9.10). We know that $\mathbf{t}^* = D^\top \theta^*$, according to the fact that $\| \mathbf{t}_{\mathcal{G}_\ell} \|_2^a - \| \mathbf{t}_{\mathcal{G}_\ell} - \mathbf{t}_{\mathcal{G}_\ell}^* \|_2^a \leq \| \mathbf{t}_{\mathcal{G}_\ell}^* \|_2^a$ and Assumption (G.4), we get

$$\text{KL}(p_0^D, \pi) \leq \ln \left(\frac{C_{\alpha,g,R}}{C_L} \right) + \alpha^a \| D^\top \theta^* \|_{a,\mathcal{G}}^a + \lambda \sum_{\ell=1}^L \ln \left\{ h(\| [D^\top \theta^*]_{\mathcal{G}_\ell} \|_2) \right\}. \quad (9.26)$$

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

$$C_{\alpha,g,R} = \int_{\Theta} \prod_{\ell=1}^L \exp(-\alpha^a \| [D^\top \theta]_{\mathcal{G}_\ell} \|_2^a) g(\| [D^\top \theta]_{\mathcal{G}_\ell} \|_2) d\theta \leq \int_{\mathbb{R}^p} \prod_{\ell=1}^L g(\| [D^\top \theta]_{\mathcal{G}_\ell} \|_2) d\theta = T.$$

Combining this with the bound of C_L^{-1} in (9.25), we obtain

$$\ln \left(\frac{C_{\alpha,g,R}}{C_L} \right) \leq r_L \alpha^a + \ln(2) \leq 1 + r_L \alpha^a. \quad (9.27)$$

Inserting (9.27) into (9.26), we get inequality (9.19). This completes the proof. \square

9.4.2. Proofs of corollaries

Proof of Corollary 5.1 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 4.1, g satisfies Assumptions (G.2) and (G.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

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

Therefore, g satisfies Assumptions (G.1)-(G.4) for $G = 1$. Owing to Remark 4.3 and Lemma 2.1, we obtain

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

We are now in position to apply Theorem 5.1 with $D = \mathbf{I}_p$ (then $q = p$), $G = 1$ (then $L = q$), $a = 1$ and $\alpha \leq 1/(3p\tau)$ to conclude. Namely, since $\tau^2 \sim (pn)^{-1}$ and $R \sim p\tau$, we get that $\Psi_{\mu_n, L, p}^D \leq 2eC_{f,\mathcal{L}}\tau^2 p \sim n^{-1}$, and

$$\Omega_{\mu_n, n, L, \lambda}^D(\theta) = \frac{4\beta}{n} \sum_{j=1}^p \ln(1 + |\theta_j|/\tau) \leq \frac{4\beta}{n} \|\theta\|_0 \ln(1 + R/\tau) \sim \frac{\|\theta\|_0 \ln(p)}{n}.$$

This completes the proof. \square

Proof of Corollary 5.2 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 4.1, g satisfies Assumptions (G.2) and (G.3). Recall that $b \in]0, 1]$. Taking $h(x) = 1 + (x/\tau)^b$ and $\lambda = c$, for all $(t, t^*) \in \mathbb{R}^G \times \mathbb{R}^G$, we have

$$\frac{g(\|t - t^*\|_2)}{g(\|t\|_2)} = \left[\frac{\tau^b + \|t\|_2^b}{\tau^b + \|t - t^*\|_2^b} \right]^c \leq \left[\frac{\tau^b + \|t - t^*\|_2^b + \|t^*\|_2^b}{\tau^b + \|t - t^*\|_2^b} \right]^c \leq \left[1 + \frac{\|t^*\|_2^b}{\tau^b + \|t - t^*\|_2^b} \right]^c \leq h(\|t^*\|_2)^\lambda.$$

Therefore, g satisfies Assumptions (G.1)-(G.4) with any $G \geq 1$. Applying Theorem 5.1, we conclude the proof. \square

Proof of Corollary 5.3 Since g satisfies Assumptions (G.2), (G.3) and D is invertible, by Remark 4.3 and Lemma 2.1, we get

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

Since $\tau^2 \sim (pn)^{-1}$ and $R \sim L\tau^a$, we get that $\Psi_{\mu_n, L, p}^D \leq 2C_{f, \mathcal{L}} \tilde{K}_{1,g}^D e\kappa^{-1} p\tau^2 \sim n^{-1}$, and

$$\Omega_{\mu_n, n, L, \lambda}^D(\theta) = \frac{c\beta}{n} \sum_{\ell=1}^L \ln \left(1 + \left[\frac{\|D^\top \theta\|_{\mathcal{G}_\ell}}{\tau} \right]^b \right) \leq \frac{c\beta}{n} \|D^\top \theta\|_{0, \mathcal{G}} \ln \left(1 + \left[\frac{R^{1/a}}{\tau} \right]^b \right) \sim \frac{\|D^\top \theta\|_{0, \mathcal{G}} \ln(L)}{n}.$$

This ends the proof. \square

9.5. Proofs of Section 6

9.5.1. Proof of Lemma 6.1

w_λ is clearly increasing, bounded from below by $w_\lambda(0)$, and continuously differentiable (in fact even C^∞) on $]0, +\infty[$.

- (i) The expression of the proximal mapping follows from [40, Lemma 7.2] provided that it is single-valued. To show the latter statement, it is sufficient to prove that the function $u : x \in [0, +\infty[\mapsto x + \gamma w_\lambda'(x) = x + \gamma c/(\tau + x) + \gamma \alpha$ has a unique minimizer occurring at 0. One can see that u admits a local maximum at $\bar{x} = -\sqrt{\gamma c} - \tau \notin [0, +\infty[$ and a local minimum at $\bar{x} = \sqrt{\gamma c} - \tau$. Thus, the problem $\min_{x \in [0, +\infty[} u(x)$ has a unique solution at 0 when $\bar{x} \leq 0$ equivalent to $\gamma \leq \tau^2/c$.
- (ii) For $\gamma \leq \tau^2/c$, it is immediate to see that the function u defined in (i) is nondecreasing on $[0, +\infty[$, and thus so is $\text{prox}_{\gamma w_\lambda}$. It then follows that $\text{prox}_{\gamma w_\lambda}(x) \geq \text{prox}_{\gamma w_\lambda}(0) = 0$ for any $x \in [0, +\infty[$. Since w_λ is increasing, the second inequality follows.
- (iii) See [40, Lemma 7.1].
- (iv) Denote the function $x \in]0, +\infty[\mapsto v(x) = \frac{x}{r} + w_\lambda'(x)$ for $r > 0$. For $x, z \geq 0$, we have

$$(v(z) - v(x))(z - x) = (z - x)^2/(r) - c \frac{(z - x)^2}{(\tau + z)(\tau + x)} \geq (z - x)^2(r^{-1} - c\tau^{-2}) \geq 0, \quad \forall r \in]0, \tau^2/c].$$

This shows that v is nondecreasing, or equivalently, that $w_\lambda(x) + \frac{1}{2r}x^2$ is convex (i.e., w_λ is semi-convex). Thus

$$\frac{\|u\|_2^2}{2r} + W_\lambda(u) = \sum_{\ell=1}^L \left(\frac{\|u_{\mathcal{G}_\ell}\|_2^2}{2r} + w_\lambda(\|u_{\mathcal{G}_\ell}\|_2) \right) = \sum_{\ell=1}^L \left\{ \left(\frac{(\cdot)^2}{2r} + w_\lambda \right) \circ \|\cdot\|_2(u_{\mathcal{G}_\ell}) \right\}.$$

Since $x \mapsto \frac{x^2}{2r} + w_\lambda(x)$ is convex and increasing, and the norm is also convex, we deduce that $\frac{\|\cdot\|_2^2}{2r} + W_\lambda$ is convex, and so is $\left(\frac{\|\cdot\|_2^2}{2r'} + W_\lambda \right) \circ D^\top$ for $r' = r/\nu$ (recall ν from (H.2)). It then follows from [40, Lemma 5.3] that $\text{prox}_{\gamma w_\lambda \circ D^\top}$ is Lipschitz continuous for any $\gamma \in]0, r]$.

- (v) [40, Lemma 3.2], which applies thanks to (i), yields the expression of $\text{prox}_{\gamma W_\lambda}$. $\nabla(\gamma W_\lambda \circ \mathbf{D}^\top)$ is Lipschitz continuous thanks to (iii), which, together with Lipschitz continuity of ∇F_β yields that of ∇V_γ . Let us turn to uniform boundedness of $\nabla(\gamma W_\lambda \circ \mathbf{D}^\top)$. For any $\boldsymbol{\theta} \in \mathbb{R}^p$, denote $\mathbf{u} = \mathbf{D}^\top \boldsymbol{\theta}$. Thus, using (i) and (iii) and that w_λ' is a decreasing function on $]0, +\infty[$, we have $\forall \boldsymbol{\theta} \in \mathbb{R}^p$,

$$\begin{aligned} \|\nabla(\gamma W_\lambda \circ \mathbf{D}^\top)(\boldsymbol{\theta})\|_2^2 &\leq \gamma^{-2} \nu \sum_{\ell=1}^L \left\{ \|\mathbf{u}_{\mathcal{G}_\ell}\|_2 - \text{prox}_{\gamma w_\lambda}(\|\mathbf{u}_{\mathcal{G}_\ell}\|_2) \right\}^2 \\ &\leq \gamma^{-2} \nu \sum_{\ell=1}^L \begin{cases} \gamma w_\lambda'(\text{prox}_{\gamma w_\lambda}(\|\mathbf{u}_{\mathcal{G}_\ell}\|_2)) & \text{if } \|\mathbf{u}_{\mathcal{G}_\ell}\|_2 > \gamma w_\lambda'(0^+) \\ \gamma w_\lambda'(0^+) & \text{otherwise} \end{cases} \\ &\leq \gamma^{-1} \nu L w_\lambda'(0^+) = \gamma^{-1} \nu L(\alpha + c/\tau). \end{aligned}$$

- (vi) Combine (i), [40, Lemma 3.2 and Lemma 6.1], (iii) and that \mathbf{M}_γ is positive definite. □

9.5.2. Proof of Theorem 6.1

Uniform geometric ergodicity We start by showing that $\mathbf{L}(t)$ in (6.5) is H -uniformly geometrically ergodic, where $H : \mathbb{R}^p \rightarrow [1, +\infty[$ is a measurable function. Let $|f|_H \stackrel{\text{def}}{=} \sup_{\boldsymbol{\theta} \in \mathbb{R}^p} \frac{|f(\boldsymbol{\theta})|}{H(\boldsymbol{\theta})}$. Recall that H -uniform geometric ergodicity requires that for all \mathbf{l}_0

$$\left| \mathbb{E} \{f(\mathbf{L}(t)) | \mathbf{L}_0 = \mathbf{l}_0\} - \int_{\mathbb{R}^p} f(\boldsymbol{\theta}) \mu_\gamma(\boldsymbol{\theta}) d\boldsymbol{\theta} \right| \leq K H(\mathbf{l}_0) \rho^t \quad (9.28)$$

for some $K < +\infty$ and $\rho \in [0, 1[$ and any function f satisfying $|f|_H \leq 1$.

Denote $\mathcal{A} \stackrel{\text{def}}{=} -\langle \nabla V_\gamma, \nabla \rangle + \Delta$ is the μ_γ symmetric natural operator. By [51, Theorem 2.2], $\mathbf{L}(t)$ is H -uniformly geometrically ergodic if H is a Lyapunov function for the generator \mathcal{A} , i.e. H is a C^2 function and there exist $\delta > 0$, $\vartheta \geq 0$ and some $R > 0$ such that for all $\boldsymbol{\theta} \in \mathbb{R}^p$

$$\mathcal{A}H(\boldsymbol{\theta}) \leq -\delta H(\boldsymbol{\theta}) + \vartheta I_{\mathbb{B}_R}(\boldsymbol{\theta}) \quad (9.29)$$

where $\mathbb{B}_R \stackrel{\text{def}}{=} \{\boldsymbol{\theta} \in \mathbb{R}^p : \|\boldsymbol{\theta}\|_2 \leq R\}$. We denote $\mathbb{B}_R^c \stackrel{\text{def}}{=} \mathbb{R}^p \setminus \mathbb{B}_R$. To prove geometric ergodicity, we use the following fact. Since F_β is differentiable, convex and e^{-F_β} is integrable, it follows from [2, Lemma 2.2] that there exist $\eta > 0$ and $r > 0$ such that for all $\boldsymbol{\theta} \in \mathbb{B}_r^c$,

$$\langle \nabla F_\beta(\boldsymbol{\theta}), \boldsymbol{\theta} \rangle \geq \eta \|\boldsymbol{\theta}\|_2. \quad (9.30)$$

Lemma 9.3. *The drift condition (9.29) is satisfied with the Lyapunov function $H(\boldsymbol{\theta}) = \exp\left(\eta/4 \sqrt{1 + \|\boldsymbol{\theta}\|_2^2}\right)$, $\delta = \eta^2/16$, $\vartheta = \eta/4H(R)\left(\eta/4 + d + 2\beta^{-1}\|\mathbf{X}\|^2 R^2 + 2\beta^{-1}\|\mathbf{X}^\top \mathbf{y}\|_2 R\right) + \eta^2/16$ and $R = \max(r, 1 + 2p/\eta)$. In turn, $\mathbf{L}(t)$ in (6.5) satisfies (9.28).*

Proof of Lemma 9.3. Elementary derivations give

$$\mathcal{A}H(\boldsymbol{\theta}) = -\eta/4 \frac{H(\boldsymbol{\theta})}{\sqrt{1 + \|\boldsymbol{\theta}\|_2^2}} \left(\langle \nabla F_\beta(\boldsymbol{\theta}), \boldsymbol{\theta} \rangle + \langle \nabla(\gamma W_\lambda)(\mathbf{D}^\top \boldsymbol{\theta}), \mathbf{D}^\top \boldsymbol{\theta} \rangle - \eta/4 \frac{\|\boldsymbol{\theta}\|_2^2}{\sqrt{1 + \|\boldsymbol{\theta}\|_2^2}} - p + \frac{\|\boldsymbol{\theta}\|_2^2}{1 + \|\boldsymbol{\theta}\|_2^2} \right).$$

Denote $\mathbf{u} = \mathbf{D}^\top \boldsymbol{\theta}$. In view of Lemma 6.1(ii), (iii) and (v), we have

$$\langle \nabla(\gamma W_\lambda)(\mathbf{u}), \mathbf{u} \rangle = \gamma^{-1} (\|\mathbf{u}\|_2^2 - \langle \text{prox}_{\gamma W_\lambda}(\mathbf{u}), \mathbf{u} \rangle) = \gamma^{-1} \sum_{\ell=1}^L \|\mathbf{u}_{\mathcal{G}_\ell}\|_2 \left\{ \|\mathbf{u}_{\mathcal{G}_\ell}\|_2 - \text{prox}_{\gamma w_\lambda}(\|\mathbf{u}_{\mathcal{G}_\ell}\|_2) \right\} \geq 0. \quad (9.31)$$

This, together with (9.30) yields the bound

$$\mathcal{A}H(\boldsymbol{\theta}) \leq -\eta/4H(\boldsymbol{\theta}) \left(\frac{\eta R - p}{1 + R} - \eta/4 \right) I_{\mathbb{B}_R^c}(\boldsymbol{\theta}) + \eta/4H(R) \left(\eta/4 + d + 2\beta^{-1}\|\mathbf{X}\|^2 R^2 + 2\beta^{-1}\|\mathbf{X}^\top \mathbf{y}\|_2 R \right) I_{\mathbb{B}_R}(\boldsymbol{\theta}),$$

whence we deduce the desired claim. □

Let us now turn to the proof of the theorem.

(i) By the triangle inequality, we have

$$\left\| \mathbf{P}_{\tilde{\mathbf{L}}_\delta}^{T,\delta,\gamma}(\mathbf{I}_0, \cdot) - \mu_n \right\|_{\text{TV}} \leq \left\| \mathbf{P}_{\tilde{\mathbf{L}}_\delta}^{T,\delta,\gamma}(\mathbf{I}_0, \cdot) - \mathbf{P}_{\mathbf{L}}^{T,\gamma}(\mathbf{I}_0, \cdot) \right\|_{\text{TV}} + \left\| \mathbf{P}_{\mathbf{L}}^{T,\gamma}(\mathbf{I}_0, \cdot) - \mu_\gamma \right\|_{\text{TV}} + \left\| \mu_\gamma - \mu_n \right\|_{\text{TV}}.$$

In view of Lemma 6.1(v) and [17, Lemma 2] and the Pinsker inequality, the first term in the right hand side goes to 0 as $\delta \rightarrow 0$. The second term converges to 0 as $T \rightarrow +\infty$ thanks to uniform geometric ergodicity (see Lemma 9.3), where we apply (9.28) with $\|f\|_\infty \leq 1$. The last term vanishes as $\gamma \rightarrow 0$ by virtue of Proposition 6.1.

(ii) In the following, C is any positive constant that does not depend on δ and T . Let $\delta \in]0, 2\beta/\|\mathbf{X}\|^2[$. Since $\mathbb{E}\{\mathbf{Z}_k\} = 0$ and $\nabla(\gamma W_\lambda) \circ \mathbf{D}^\top$ is uniformly bounded, we have

$$\begin{aligned} \left\| \mathbb{E}\{\mathbf{L}_{k+1}\} \right\|_2 &\leq \left\| \mathbb{E}\left\{ \mathbf{L}_k - \delta/2 \nabla F_\beta(\mathbf{L}_k) \right\} \right\|_2 + C\delta \\ &\leq \mathbb{E}\left\{ \left\| (\mathbf{I}_p - \delta/2 \nabla F_\beta)(\mathbf{L}_k) - (\mathbf{I}_p - \delta/2 \nabla F_\beta)(0) \right\|_2 \right\} + (C + \left\| \nabla F_\beta(0) \right\|_2/2)\delta \\ &\leq \mathbb{E}\{\mathbf{L}_k\} + (C + \left\| \mathbf{X}^\top \mathbf{y} \right\|_2/2)\delta \leq (C + \left\| \mathbf{X}^\top \mathbf{y} \right\|_2/2)T, \quad \forall k[0, \lfloor T/\delta \rfloor], \end{aligned}$$

where we used the fact that F_β is a differentiable convex function whose gradient is $2\|\mathbf{X}\|^2/\beta$ -Lipschitz, and thus, $\mathbf{I}_p - \delta \nabla F_\beta$ is non-expansive for the prescribed choice of δ [3]. In addition, by independence of \mathbf{Z}_k from \mathbf{L}_k and \mathbf{Y} , and in view of (9.31) and Lipschitz continuity of ∇F_β , we have

$$\begin{aligned} \mathbb{E}\left\{ \left\| \mathbf{L}_{k+1} \right\|_2^2 \right\} &\leq \mathbb{E}\left\{ \left\| \mathbf{L}_k - \delta/2 \nabla V_\gamma(\mathbf{L}_k) \right\|_2^2 \right\} + \delta p \\ &\leq \mathbb{E}\left\{ \left\| \mathbf{L}_k - \delta/2 \nabla F_\beta(\mathbf{L}_k) + \delta/2 \nabla F_\beta(0) \right\|_2^2 - \delta \left\langle \mathbf{L}_k - \delta/2 (\nabla F_\beta(\mathbf{L}_k) - \nabla F_\beta(0)), \nabla(\gamma W_\lambda \circ \mathbf{D}^\top)(\mathbf{L}_k) \right\rangle \right\} \\ &\quad + C\delta^2 + \delta p \\ &\leq \mathbb{E}\left\{ \left\| \mathbf{L}_k \right\|_2^2 - \delta \left\langle \mathbf{L}_k - \delta/2 (\nabla F_\beta(\mathbf{L}_k) - \nabla F_\beta(0)), \nabla(\gamma W_\lambda \circ \mathbf{D}^\top)(\mathbf{L}_k) \right\rangle \right\} + C\delta^2 + \delta p \\ &\leq \mathbb{E}\left\{ \left\| \mathbf{L}_k \right\|_2^2 \right\} + C\delta^2 \left\| \mathbb{E}\{\mathbf{L}_k\} \right\|_2 + C\delta^2 + \delta p \\ &\leq \mathbb{E}\left\{ \left\| \mathbf{L}_k \right\|_2^2 \right\} + C\delta^2 \left\| \mathbb{E}\{\mathbf{L}_k\} \right\|_2 + C\delta^2 + \delta p \\ &\leq \mathbb{E}\left\{ \left\| \mathbf{L}_k \right\|_2^2 \right\} + C\delta^2 T + C\delta^2 + \delta p \leq C\delta T^2 + C\delta T + pT = O(T^2), \quad \forall k[0, \lfloor T/\delta \rfloor]. \end{aligned} \quad (9.32)$$

For $T > 0$ and $R > 0$, denote $\tilde{\mathbf{L}}_\delta^T \stackrel{\text{def}}{=} 1/T \int_0^T \tilde{\mathbf{L}}_\delta(t) dt$, $\tilde{\mathbf{L}}_\delta^{T,R} \stackrel{\text{def}}{=} 1/T \int_0^T \tilde{\mathbf{L}}_\delta(t) I_{\mathbb{B}_R}(\tilde{\mathbf{L}}_\delta(t)) dt$, $\mathbf{L}^{T,R} \stackrel{\text{def}}{=} 1/T \int_0^T \mathbf{L}(t) I_{\mathbb{B}_R}(\mathbf{L}(t)) dt$ and $\widehat{\boldsymbol{\theta}}_\gamma^R \stackrel{\text{def}}{=} \int_{\mathbb{B}_R} \boldsymbol{\theta} \mu_\gamma(\boldsymbol{\theta}) d\boldsymbol{\theta}$. The triangle and Jensen inequalities yield

$$\begin{aligned} \mathbb{E}\left\{ \left\| \tilde{\mathbf{L}}_{\delta,T,\gamma} - \widehat{\boldsymbol{\theta}}_\gamma^R \right\|_2 \right\} &\leq \mathbb{E}\left\{ \left\| \tilde{\mathbf{L}}_{\delta,T,\gamma} - \tilde{\mathbf{L}}_\delta^T \right\|_2 \right\} + \mathbb{E}\left\{ \left\| \tilde{\mathbf{L}}_\delta^{T,R} - \mathbf{L}^{T,R} \right\|_2 \right\} + \mathbb{E}\left\{ \left\| \mathbf{L}^{T,R} - \widehat{\boldsymbol{\theta}}_\gamma^R \right\|_2 \right\} \\ &\quad + \frac{1}{T} \int_0^T \mathbb{E}\left\{ \left\| \tilde{\mathbf{L}}_\delta(t) \right\|_2 I_{\mathbb{B}_R^c}(\tilde{\mathbf{L}}_\delta(t)) \right\} dt + \int_{\mathbb{B}_R^c} \left\| \boldsymbol{\theta} \right\|_2 \mu_\gamma(\boldsymbol{\theta}) d\boldsymbol{\theta}. \end{aligned}$$

For the last two terms, we have the bounds

$$\begin{aligned} \frac{1}{T} \int_0^T \mathbb{E}\left\{ \left\| \tilde{\mathbf{L}}_\delta(t) \right\|_2 I_{\mathbb{B}_R^c}(\tilde{\mathbf{L}}_\delta(t)) \right\} dt &\leq R^{-1} \frac{1}{T} \int_0^T \mathbb{E}\left\{ \left\| \tilde{\mathbf{L}}_\delta(t) \right\|_2^2 \right\} dt \stackrel{(9.32)}{\leq} CT^2/R \quad \text{and} \\ \int_{\mathbb{B}_R^c} \left\| \boldsymbol{\theta} \right\|_2 \mu_\gamma(\boldsymbol{\theta}) d\boldsymbol{\theta} &\leq R^{-1} \int_{\mathbb{R}^p} \left\| \boldsymbol{\theta} \right\|_2^2 \mu_\gamma(\boldsymbol{\theta}) d\boldsymbol{\theta} \stackrel{(G.3)}{\leq} C/R. \end{aligned}$$

Choosing, e.g., $R = 1/T^3$, these terms converge to 0 as $T \rightarrow \infty$. Using (9.32) and arguing as in [23, Step 1, Proposition 2], we have

$$\mathbb{E}\left\{ \left\| \tilde{\mathbf{L}}_{\delta,T,\gamma} - \tilde{\mathbf{L}}_\delta^T \right\|_2 \right\} \leq C\delta(1 + \delta T^2) \rightarrow 0 \text{ as } \delta \rightarrow 0.$$

Using Girsanov formula and Pinsker inequality, as in [23, Step 2, Proposition 2], the distribution of $\{\tilde{L}_\delta(t)\}_{t \in [0, T]}$ converges to that of $\{L(t)\}_{t \in [0, T]}$ in total variation as $\delta \rightarrow 0$. In turn,

$$\mathbb{E} \left\{ \left\| \tilde{L}_\delta^{T, R} - L^{T, R} \right\|_2 \right\} \rightarrow 0 \text{ as } \delta \rightarrow 0.$$

Using Lemma 9.3 with $f(\theta) = \theta_i I_{\mathbb{B}_R}(\theta)$, and arguing as in [23, Step 4, Proposition 2], we have

$$\mathbb{E} \left\{ \left\| L^{T, R} - \tilde{\theta}_y^R \right\|_2 \right\} \leq CT^{-1/2} \rightarrow 0 \text{ as } T \rightarrow \infty.$$

This completes the proof.

Acknowledgements. 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] D. Bakry, F. Barthe, P. Cattiaux, and A. Guillin. A simple proof of the Poincaré inequality for a large class of probability measures. *Electronic Communications in Probability*, 13:60–66, 2008.
- [3] H. H. Bauschke and P. L. Combettes. *Convex analysis and monotone operator theory in Hilbert spaces*. Springer, 2011.
- [4] G. Biau. Analysis of a random forests model. *J. Mach. Learn. Res.*, 13(1):1063–1095, Apr. 2012.
- [5] 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.
- [6] G. Biau, L. Devroye, and G. Lugosi. Consistency of random forests and other averaging classifiers. *J. Mach. Learn. Res.*, 9:2015–2033, June 2008.
- [7] L. Breiman. Bagging predictors. *Mach. Learn.*, 24(2):123–140, Aug. 1996.
- [8] L. Breiman. Random forests. *Mach. Learn.*, 45(1):5–32, Oct. 2001.
- [9] P. Bühlmann and S. Van De Geer. *Statistics for high-dimensional data: methods, theory and applications*. Springer, 2011.
- [10] 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.
- [11] E. J. Candès. Ridgelets: Estimating with Ridge Functions. *Annals of Statistics*, 31, 1999. 1561–1599.
- [12] L. Chaari, J.-Y. Tournet, C. Chaux, and H. Batatia. A hamiltonian monte carlo method for non-smooth energy sampling. Technical Report arXiv:1401.3988, , 2014.
- [13] X. Chen, Q. Lin, S. Kim, J. G. Carbonell, and E. P. Xing. An efficient proximal-gradient method for general structured sparse learning. Preprint arXiv:1005.4717, 2010.
- [14] C. Chesneau, M. Fadili, and J.-L. Starck. Stein block thresholding for image denoising. *Applied and Computational Harmonic Analysis*, 28(1):67–88, 2010.
- [15] R. Coifman and D. Donoho. Translation invariant de-noising. In A. Antoniadis and G. Oppenheim, editors, *Wavelets and Statistics*. Springer-Verlag, 1995. 125–150.
- [16] D. Dai, P. Rigollet, and T. Zhang. Deviation optimal learning using greedy q -aggregation. *Ann. Statist.*, 40(3):1878–1905, 06 2012.
- [17] A. Dalalyan. Theoretical guarantees for approximate sampling from smooth and log-concave densities. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(3):651–676, 2016.
- [18] 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.
- [19] 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.
- [20] A. S. Dalalyan and J. Salmon. Sharp oracle inequalities for aggregation of affine estimators. *Ann. Statist.*, 40(4):2327–2355, 08 2012.
- [21] 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.
- [22] A. S. Dalalyan and A. B. Tsybakov. Mirror averaging with sparsity priors. *Bernoulli*, 18(3):914–944, 08 2012.
- [23] 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.
- [24] 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.
- [25] D. L. Donoho and I. M. Johnstone. Ideal spatial adaptation by wavelet shrinkage. *Biometrika*, 81:425–455, 1994.
- [26] D. L. Donoho and I. M. Johnstone. Adapting to unknown smoothness via wavelet shrinkage. *Journal of the American Statistical Association*, pages 1200–1224, 1995.
- [27] D. L. Donoho, I. M. Johnstone, G. Kerkycharian, and D. Picard. Wavelet shrinkage: asymptopia. *Journal of the Royal Statistical Society, Ser. B*, pages 371–394, 1995.
- [28] A. Durmus and E. Moulines. Non-asymptotic convergence analysis for the Unadjusted Langevin Algorithm. Preprint hal-01176132, July 2015.

- [29] A. Durmus, E. Moulines, and M. Pereyra. Sampling from convex non continuously differentiable functions, when Moreau meets Langevin. Preprint hal-01267115, Feb. 2016.
- [30] K. Fang, S. Kotz, and K. Ng. *Symmetric multivariate and related distributions*. Monographs on statistics and applied probability. Chapman and Hall, 1990.
- [31] Y. Freund. Boosting a weak learning algorithm by majority. *Information and Computation*, 121(2):256 – 285, 1995.
- [32] R. Genuer. *Random Forests: elements of theory, variable selection and applications*. Theses, Université Paris Sud - Paris XI, Nov. 2010.
- [33] I. S. Gradshteyn and I. M. Ryzhik. *Table of Integrals, Series, and Products*. Academic Press, London, 4th edition, 1965.
- [34] J. Huang, J. L. Horowitz, and F. Wei. Variable selection in nonparametric additive models. *Ann. Statist.*, 38(4):2282–2313, 08 2010.
- [35] J. Kivinen and M. K. Warmuth. Exponentiated gradient versus gradient descent for linear predictors. *Inf. Comput.*, 132(1):1–63, Jan. 1997.
- [36] V. Koltchinskii and M. Yuan. Sparse recovery in large ensembles of kernel machines, 2008.
- [37] G. Lecué. Simultaneous adaptation to the margin and to complexity in classification. *Ann. Statist.*, 35(4):1698–1721, 08 2007.
- [38] G. Leung and A. R. Barron. Information theory and mixing least-squares regressions. *IEEE Trans. Inform. Theory*, 52(8):3396–3410, 2006.
- [39] N. Littlestone and M. K. Warmuth. The weighted majority algorithm. *Inf. Comput.*, 108(2):212–261, Feb. 1994.
- [40] T. D. Luu, J. Fadili, and C. Chesneau. Sampling from non-smooth distribution through Langevin diffusion. Technical report, hal-01492056, May 2017.
- [41] J. Maly. Lectures on change of variables in integral. Preprint 305, Department of Mathematics, University of Helsinki, 2001.
- [42] L. Meier, S. Van De Geer, and P. Bühlmann. 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] P. Müller and B. Vidakovic, editors. *Bayesian Inference in Wavelet-Based Models*, volume 141 of *Lecture Notes in Statistics*. Springer-Verlag, New York, 1999.
- [45] A. Nemirovski. Topics in non-parametric statistics, 2000.
- [46] M. Pereyra. Proximal markov chain monte carlo algorithms. *Statistics and Computing*, 26(4):745–760, 2016.
- [47] G. Peyré, J. Fadili, and C. Chesneau. Group sparsity with overlapping partition functions. In *EUSIPCO*, Barcelona, Spain, Aug. 2011.
- [48] 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.
- [49] P. Rigollet and A. Tsybakov. Linear and convex aggregation of density estimators. *Mathematical Methods of Statistics*, 16(3):260–280, 2007.
- [50] P. Rigollet and A. B. Tsybakov. Sparse estimation by exponential weighting. *Statist. Sci.*, 27(4):558–575, 11 2012.
- [51] G. O. Roberts and R. L. Tweedie. Exponential Convergence of Langevin Distributions and Their Discrete Approximations. *Bernoulli*, 2(4):341–363, 1996.
- [52] R. T. Rockafellar and R. Wets. *Variational analysis*, volume 317. Springer Verlag, 1998.
- [53] L. I. Rudin, S. Osher, and E. Fatemi. Nonlinear total variation based noise removal algorithms. *Phys. D*, 60(1-4):259–268, Nov. 1992.
- [54] R. E. Schapire. The strength of weak learnability. *Mach. Learn.*, 5(2):197–227, July 1990.
- [55] 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.
- [56] 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.
- [57] M. Xuerong. *Stochastic differential equations and applications*. Woodhead Publishing, 2007.
- [58] Y. Yang. Aggregating regression procedures to improve performance. *Bernoulli*, 10(1):25–47, 2004.
- [59] 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.