

n° 2014-26
**Concentration of Quadratic
Forms and Aggregation of
Affine Estimators**

P. BELLEC¹

Les documents de travail ne reflètent pas la position du CREST et n'engagent que leurs auteurs.
Working papers do not reflect the position of CREST but only the views of the authors.

¹ CREST-ENSAE, 3 Avenue Pierre Larousse, 92245 Malakoff Cedex, France and CMAP-Ecole Polytechnique, 91120 Palaiseau, France.

Concentration of quadratic forms and aggregation of affine estimators

Pierre Bellec ^{*†‡}

October 2, 2014

Abstract

This paper deals with aggregation of estimators in the context of regression with fixed design, with heteroscedastic and subgaussian noise. We relate the task of aggregating a finite family of affine estimators to the concentration of quadratic forms of the noise vector, and we derive sharp oracle inequalities in deviation for model selection type aggregation of affine estimators when the noise is subgaussian. Explicit numerical constants are given for Gaussian noise. Then we present a new concentration result that is sharper than the Hanson-Wright inequality under the Bernstein condition on the noise. This allows us to improve the sharp oracle inequality obtained in the subgaussian case. Finally, we show that up to numerical constants, the optimal sparsity oracle inequality previously obtained for Gaussian noise holds in the subgaussian case. The exact knowledge of the variance of the noise is not needed to construct the estimator that satisfies the sparsity oracle inequality.

1 Introduction

We study the problem of recovering an unknown vector $\mathbf{f} = (f_1, \dots, f_n)^T \in \mathbf{R}^n$ from noisy observations

$$Y_i = f_i + \xi_i, \quad i = 1, \dots, n, \quad (1.1)$$

where the noise random variables ξ_1, \dots, ξ_n are zero mean, subgaussian random variables. We measure the quality of estimation of the unknown vector \mathbf{f} with the squared euclidean norm in \mathbf{R}^n :

$$\|\mathbf{f} - \hat{\mathbf{f}}\|_2^2,$$

for any estimator $\hat{\mathbf{f}}$ of \mathbf{f} . When the noise random variables are normal, this is the Gaussian sequence model, which has been extensively studied [20]. Several estimators have

^{*}CREST-ENSAE, 3 avenue Pierre Larousse, 92245 Malakoff Cedex, France

[†]CMAP, Ecole Polytechnique, 91120 Palaiseau, France

[‡]This work was supported by the grant Investissements d'Avenir (ANR-11-IDEX-0003/Labex Ecodec/ANR-11-LABX-0047).

been proposed to recover the unknown vector \mathbf{f} from the observations: the Ordinary Least Squares, the Ridge regressors, the Stein estimator and the procedures based on shrinkage, to name a few. Several of these estimators depend on a parameter that must be chosen carefully to obtain satisfying error bounds. These available estimators have different strengths and weaknesses in different scenarios, so it is important to be able to mimic the best among a given family of estimators, without any assumption on the unknown regression vector \mathbf{f} . The problem of mimicking the best estimator in a given finite set is the problem of model-selection type aggregation, which was introduced in [23, 29]. More precisely, let $\hat{\mu}_1, \dots, \hat{\mu}_M$ be M estimators of \mathbf{f} based on the data Y_1, \dots, Y_n . The goal is to construct a new estimator or aggregate \hat{f} with the same data Y_1, \dots, Y_n , which satisfies with probability greater than $1 - \epsilon$:

$$\left\| \hat{f} - \mathbf{f} \right\|_2^2 \leq \min_{j=1, \dots, M} \left\| \hat{\mu}_j - \mathbf{f} \right\|_2^2 + \delta_{n, M}(\epsilon),$$

where $\delta_{n, M}(\cdot)$ is a function of ϵ that should be small. The above inequality is called a sharp oracle inequality. Here, *sharp* means that the coefficient of the oracle risk $\min_{j=1, \dots, M} \left\| \hat{\mu}_j - \mathbf{f} \right\|_2^2$ is 1, which is essential to derive minimax optimality results.

A first approach to mimic the best estimators in a given family is to use independence by assuming that the estimators $\hat{\mu}_1, \dots, \hat{\mu}_n$ are independent of the observations Y_1, \dots, Y_n used for the aggregation step. For example, assume that two independent samples (Y_1, \dots, Y_n) and (Y'_1, \dots, Y'_n) are available, with Y_i and Y'_i independent and identically distributed for all $i = 1, \dots, n$. Then one can use the sample Y_1, \dots, Y_n to construct the estimators $\hat{\mu}_1, \dots, \hat{\mu}_M$ and use the independent sample Y'_1, \dots, Y'_n to aggregate them. For the aggregation step, conditionally on Y_1, \dots, Y_n , the estimators $\hat{\mu}_1, \dots, \hat{\mu}_M$ can be considered deterministic, thanks to independence. It is possible to obtain such independent samples when the noise is Gaussian and the variance is known, with *sample cloning* [28, Lemma 2.1], at the cost of a factor 2 in the variance of the observations. However, this technique is specific to the Gaussian case and cannot be used when the noise is only assumed to be subgaussian as in the present paper.

Among the procedures available to estimate \mathbf{f} , several are linear in the observations Y_1, \dots, Y_n . It is the case for example of the Least Squares and the Ridge regressors, whereas the shrinkage estimators and the Stein estimator are non-linear functions of the observations. A description of the estimators that are linear or affine in the observations is given in [11, Section 1.2], [1] and references therein. This linear behavior of the estimators $\hat{\mu}_1, \dots, \hat{\mu}_M$ makes it possible to explicitly treat the dependence between the estimators $\hat{\mu}_1, \dots, \hat{\mu}_M$ and the data Y_1, \dots, Y_n used to aggregate them. Leung and Barron [22] studied the problem of aggregation of projection estimators, and derived sharp oracle inequalities in expectation with a procedure based on exponential weights. Then, Dalalyan and Salmon [11] and Dai et al. [9] gave insights on how to construct an aggregate to mimic the best candidate among a set of affine estimators. Here we also consider affine estimators. Let $\mathbf{y} = (Y_1, \dots, Y_n)^T$ be the vector of observations. An affine estimator is of the form $\hat{\mu}_j = A_j \mathbf{y} + b_j$ for a deterministic matrix A_j of size $n \times n$ and a deterministic vector $b_j \in \mathbf{R}^n$.

We consider in Section 3 that the variances of the noise random variables ξ_1, \dots, ξ_n are known and in Section 4 that an upper bound on the subgaussian norm of the noise vector is known. We refer the reader to [16] and the survey [17] for the problem of

estimating the unknown vector \mathbf{f} when the variance of the noise is unknown, which is outside of the scope of the present paper.

As in the papers [11, 9], we consider the problem of aggregation of M affine estimators with a prior probability distribution π_1, \dots, π_M on the finite set of indices $\{1, \dots, M\}$. Prior weights is a common ingredient in deriving sharp oracle inequalities for model-selection type aggregation [12, 10, 21, 4]. An example of such an oracle inequality is (1.2) below. The use of sparsity-inducing prior weights is crucial to prove sparsity oracle inequalities via sparsity pattern aggregation [25, 24, 9, 28]. When the noise is Gaussian with variance σ^2 , the following sparsity oracle inequality was shown in [9] for an estimator $\hat{\mu}$ and a design matrix \mathbb{X} with p columns: with probability greater than $1 - 2 \exp(-x)$,

$$\|\hat{\mu} - \mathbf{f}\|_2^2 \leq \min_{\theta \in \mathbf{R}^p} \left(\|\mathbb{X}\theta - \mathbf{f}\|_2^2 + c \sigma^2 |\theta|_0 \log \left(\frac{ep}{1 \vee |\theta|_0} \right) \right) + c' \sigma^2 x.$$

In the previous display, $c, c' > 0$ are absolute constants and $|\theta|_0$ denotes the number of non-zero coefficients of θ . A similar result in expectation was shown in [25, 28], also with the assumption that the noise random variables are normal. In Section 4, we propose an estimator that achieves a similar sparsity oracle inequality in deviation, but we only assume that the noise vector is subgaussian. It extends the previous results [25, 24, 9, 28] to the subgaussian setting.

The papers [11, 9] derived different procedures that satisfy sharp oracle inequalities for the problem of aggregation of affine estimators when the noise random variable are Gaussian. Dalalyan and Salmon [11] proposed an estimator $\hat{\mu}^{EW}$ based on exponential weights, for which a sharp oracle inequality holds in expectation:

$$\mathbb{E} \|\mathbf{f} - \hat{\mu}^{EW}\|_2^2 \leq \min_{j=1, \dots, M} \left(\mathbb{E} \|\hat{\mu}_j - \mathbf{f}\|_2^2 + \beta \log \frac{1}{\pi_j} \right), \quad (1.2)$$

where β is a constant proportional to the largest variance of the noise random variables. This oracle inequality in expectation holds for $\hat{\mu}^{EW}$ under a commutativity assumption on the matrices A_j , which is enough to apply this oracle inequality to orthogonal projections on a set of coordinates. In the case where the matrices A_j are not symmetric, [11] achieved a similar oracle inequality by symmetrizing the affine estimators before the aggregation step, which suggests that the symmetry assumption can be relaxed. Although the estimator $\hat{\mu}^{EW}$ achieves this inequality in expectation, it was shown in [10] that this procedure cannot achieve a similar result in deviation, with an unavoidable error term of order \sqrt{n} . In Dai et al. [9], a sharp oracle inequality in deviation is derived for an estimator $\hat{\mu}^Q$ based on Q -aggregation [10]. Namely, the estimator $\hat{\mu}^Q$ satisfies with probability greater than $1 - \delta$:

$$\|\mathbf{f} - \hat{\mu}^Q\|_2^2 \leq \min_{j=1, \dots, M} \left(\|\hat{\mu}_j - \mathbf{f}\|_2^2 + 4\sigma^2 \text{Tr}(A_j) + \beta \log \frac{1}{\pi_j} \right) + \beta \log \frac{1}{\delta}, \quad (1.3)$$

where β is a constant and the noise random variables are i.i.d. with variance σ^2 . This bound shows that it is possible to achieve oracle inequalities in deviation in the context of aggregation of affine estimators. However the extra term $4\sigma^2 \text{Tr}(A_j)$ may be large in common situation where the trace of some matrices A_j is large. For example, if one

aggregates the estimators $\hat{\mu}_1 = \lambda_1 \mathbf{y}, \dots, \hat{\mu}_M = \lambda_M \mathbf{y}$, for some positive real numbers $\lambda_1, \dots, \lambda_M$ with the uniform prior $\pi_j = 1/M$ for all $j = 1, \dots, M$, then the remainder term $4\sigma^2 \text{Tr}(A_j)$ in the above oracle inequality is of order $\sigma^2 n \lambda_j$ for each $j = 1, \dots, M$, which is large relatively to the optimal rate $\sigma^2 \log M$. This term $4\sigma^2 \text{Tr}(A_j)$ makes the previous oracle inequality suitable only for scenarios where the matrices A_j have small trace.

The contributions of the present paper are the following:

- We propose an estimator that satisfies a sharp oracle inequality in deviation without the extra term proportional to $\sigma^2 \text{Tr}(A_j)$, under three different assumptions on the noise. This is our main result and it is given in Theorem 3.1. Under the three Assumptions 3.1, 3.2 and 3.3, our estimator is suitable for situations involving matrices A_j with large trace, and it recovers the optimal rate proportional to $\log M$ when the uniform prior is used. Assumption 3.1 deals with heteroscedastic Gaussian noise and then explicit absolute constants are provided for the sharp oracle inequality. Under Assumption 3.2, the noise random variables are independent and subgaussian, and the multiplicative constant β may be arbitrarily large for noise random variables with pathologically small variance. Assumption 3.3 is slightly stronger than Assumption 3.2, which prevents the variance from being too small relatively to its subgaussian norm, and under this third assumption we can control the value of β . In earlier results [11, 9], only Gaussian noise was considered.
- In order to prove Theorem 3.1 under Assumption 3.3, we derive a new concentration result for quadratic forms of independent random variables which is given in Theorem 3.2. It is sharper than the Hanson-Wright inequality under Assumption 3.3.
- The assumptions on the matrices A_1, \dots, A_M are relaxed. In particular, they can be non-symmetric and have negative eigenvalues.
- Using sparsity pattern aggregation, we derive a sparsity oracle inequality in deviation when the noise vector is subgaussian, without assuming independence of the noise components. Theorem 4.1 recovers up to absolute constants the sparsity oracle inequality obtained when the noise is Gaussian [25, 24, 9].

The paper is organized as follows. In Section 2 we define the notation used throughout the paper. Section 3 defines an estimator and shows that it achieves sharp oracle inequalities in deviation for aggregation of affine estimators under three different assumptions on the noise. In Section 4, we derive a sparsity oracle inequality when the noise vector is subgaussian. The concentration inequalities used in the paper are given in Appendix A and the proofs are given in Appendix B.

2 Notation

We study an aggregation problem for the regression model with fixed design and heteroscedastic subgaussian noise. A random variable X is said to be subgaussian if and only if the quantity

$$\|X\|_{\psi_2} = \sup_{p \geq 1} p^{-1/2} (\mathbb{E}|X|^p)^{1/p}$$

is finite. Several other definitions are used in the literature, see [30, Section 5.2.3] for a review of their equivalence.

Let $(f_1, \dots, f_n)^T \in \mathbf{R}^n$ be an unknown regression vector. We observe n random variables (1.1) where ξ_1, \dots, ξ_n are subgaussian random variables, with $\mathbb{E}[\xi_i] = 0$ and $\mathbb{E}[\xi_i^2] = \sigma_i^2$. The model is heteroscedastic, which means that the random variables ξ_1, \dots, ξ_n may have different variances. It can be rewritten in the vector form $\mathbf{y} = \mathbf{f} + \boldsymbol{\xi}$ where $\mathbf{y} = (Y_1, \dots, Y_n)^T$, $\mathbf{f} = (f_1, \dots, f_n)^T$ and $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n)^T$.

For any estimator \hat{f}_n of \mathbf{f} , we measure the quality of estimation of \mathbf{f} with the loss $\|\mathbf{f} - \hat{f}_n\|_2^2$ where $\|\cdot\|_2$ denotes the Euclidean norm in \mathbf{R}^n . Let $M \geq 2$. As in [11, 9], we consider M affine estimators of the form

$$\hat{\mu}_j = A_j \mathbf{y} + b_j, \quad j = 1, \dots, M.$$

The matrices A_1, \dots, A_M and the vectors $b_1, \dots, b_M \in \mathbf{R}^n$ are deterministic. Define the simplex in \mathbf{R}^M :

$$\Lambda^M = \left\{ \theta \in \mathbf{R}^M, \sum_{j=1}^M \theta_j = 1, \forall j = 1 \dots M, \theta_j \geq 0 \right\}$$

and for any $\theta \in \Lambda^M$, let $\hat{\mu}_\theta = \sum_{j=1}^M \theta_j \hat{\mu}_j$. Let e_1, \dots, e_M be the vectors of the canonical basis in \mathbf{R}^M . Then $\hat{\mu}_j = \hat{\mu}_{e_j}$ for all $j = 1, \dots, M$.

Finally, for any $n \times n$ real matrix $A = (a_{i,j})_{i,j=1,\dots,n}$, define the operator norm of A , the Hilbert-Schmidt (or Frobenius) norm of A and the nuclear norm of A respectively by:

$$\|A\|_2 = \sup_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2}, \quad \|A\|_{\text{HS}} = \sqrt{\sum_{i,j=1,\dots,n} a_{i,j}^2}, \quad \|A\|_1 = \text{Tr}(\sqrt{A^T A}). \quad (2.1)$$

3 Model-selection type oracle inequalities

3.1 The proposed estimator

For any $\theta \in \Lambda^M$ define

$$\begin{aligned} \hat{H}_n(\theta) &= \|\hat{\mu}_\theta\|_2^2 - 2\mathbf{y}^T \hat{\mu}_\theta + 2 \sum_{j=1}^M \theta_j \text{Tr}(D_\sigma A_j D_\sigma) \\ &\quad + \frac{1}{2} \widehat{\text{pen}}(\theta) + \beta \sum_{j=1}^M \theta_j \log \frac{1}{\pi_j}, \end{aligned} \quad (3.1)$$

where $\beta > 0$ is a constant, $D_\sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$ and

$$\widehat{\text{pen}}(\theta) = \sum_{j=1}^M \theta_j \|\hat{\mu}_\theta - \hat{\mu}_j\|_2^2. \quad (3.2)$$

We consider the estimator $\hat{\mu}_{\hat{\theta}}$ where

$$\hat{\theta} \in \underset{\theta \in \Lambda^M}{\operatorname{argmin}} \hat{H}_n(\theta). \quad (3.3)$$

When θ is fixed and deterministic, the term

$$\|\hat{\mu}_\theta\|_2^2 - 2\mathbf{y}^T \hat{\mu}_\theta + 2 \sum_{j=1}^M \theta_j \operatorname{Tr}(D_\sigma A_j D_\sigma) \quad (3.4)$$

in the definition of \hat{H}_n is an unbiased estimate of the quantity

$$\|\hat{\mu}_\theta\|_2^2 - 2\mathbf{f}^T \hat{\mu}_\theta = \|\hat{\mu}_\theta - \mathbf{f}\|_2^2 - \|\mathbf{f}\|_2^2, \quad (3.5)$$

which is the quantity of interest $\|\hat{\mu}_\theta - \mathbf{f}\|_2^2$ up to the additive constant $\|\mathbf{f}\|_2^2$. The term involving the trace of the matrices $D_\sigma A_j D_\sigma$ comes from the quadratic term in $\boldsymbol{\xi}$:

$$\sum_{j=1}^M \theta_j \operatorname{Tr}(D_\sigma A_j D_\sigma) = \mathbb{E}[\sum_{j=1}^M \theta_j \boldsymbol{\xi}^T A_j \boldsymbol{\xi}] = \mathbb{E}[\boldsymbol{\xi}^T \hat{\mu}_\theta].$$

The estimators from [22, 11, 9] are all obtained with an unbiased estimate of the quantity (3.5), so the term (3.4) comes as no surprise in the definition of \hat{H}_n .

The penalty (3.2) is borrowed from the Q -aggregation procedure, which is a powerful tool to derive sharp oracle inequalities in deviation when the loss is strongly convex [10, 21, 4]. Since the estimators $\hat{\mu}_1, \dots, \hat{\mu}_M$ depend on the data, the penalty (3.2) is data-driven, which is not the case when the estimators to aggregate are deterministic vectors as in [10]. In order to give some geometric insights on the penalty (3.2), let $c \in \mathbf{R}^n$ satisfies the M linear equations $2c^T \hat{\mu}_j = \|\hat{\mu}_j\|_2^2$ and assume only in the rest of this paragraph that c is well defined, even though this assumption cannot be fulfilled for $M > n$. Then

$$\widehat{\operatorname{pen}}(\theta) = \sum_{j=1}^M \theta_j \|\hat{\mu}_j\|_2^2 - \|\hat{\mu}_\theta\|_2^2 = 2c^T \hat{\mu}_\theta - \|\hat{\mu}_\theta\|_2^2 = \|c\|_2^2 - \|\hat{\mu}_\theta - c\|_2^2. \quad (3.6)$$

Assume also only in this paragraph that the function $\theta \rightarrow \hat{\mu}_\theta$ is bijective from the simplex Λ^M to the convex hull of $\{\hat{\mu}_1, \dots, \hat{\mu}_M\}$. Then we can write $\widehat{\operatorname{pen}}(\theta) = g(\hat{\mu}_\theta)$ for some function g defined on the convex hull of $\{\hat{\mu}_1, \dots, \hat{\mu}_M\}$. Equation (3.6) shows that the level sets of the function g are euclidean balls centered at c . The function g is non-negative, it is minimal at the extreme points $\hat{\mu}_1, \dots, \hat{\mu}_M$ since $g(\hat{\mu}_j) = 0$ for all $j = 1, \dots, M$ and g is maximal at the projection of c on the convex hull of $\{\hat{\mu}_1, \dots, \hat{\mu}_M\}$. Intuitively, the penalty (3.2) pushes θ away from the center of the simplex towards the vertices. Thus, the level sets of the function $\theta \rightarrow \widehat{\operatorname{pen}}(\theta)$ in \mathbf{R}^M are ellipsoids centered at θ_c , where θ_c is the unique point in \mathbf{R}^M such that $\hat{\mu}_{\theta_c} = c$. If $M > n$ or if the vector c is not well defined, the level sets of $\widehat{\operatorname{pen}}(\cdot)$ are more intricate and cannot be described as simply.

Finally, the term

$$\beta \sum_{j=1}^M \theta_j \log \frac{1}{\pi_j} \quad (3.7)$$

allows to weight the candidates $\hat{\mu}_1, \dots, \hat{\mu}_M$ with the prior probability distribution $(\pi_j)_{j=1, \dots, M}$ based on some prior knowledge about the estimators $\hat{\mu}_1, \dots, \hat{\mu}_M$. That prior probability distribution $(\pi_j)_{j=1, \dots, M}$ is deterministic and cannot depend on the data Y_1, \dots, Y_n . For example, if the estimators are projection estimators, one can set prior weights that decrease with the rank of the projections [24], we use this strategy in Section 4. The same term is used in [21] whereas [9] uses the Kullback-Leibler divergence of θ from π . It is shown in [10] that for aggregation of deterministic vectors, one may use a quantity of the form $\beta \sum_{j=1}^M \theta_j \log(\rho(\theta_j)/\pi_j)$ where $\rho(\cdot)$ satisfies $\rho(t) \geq t$ and $t \rightarrow t \log(\rho(t))$ is convex. This suggests that we could use the Kullback-Leibler divergence of θ from π instead of (3.7), but in their current form, our proofs only hold with the “linear entropy” (3.7).

Finally, notice that the function \hat{H}_n is convex, as it has the form $\hat{H}_n(\theta) = \frac{1}{2} \|\hat{\mu}_\theta\|_2^2 + \text{lin}(\theta)$ where $\text{lin}(\cdot)$ is a linear function. This can be seen using (B.2) with $g = 0$. Thus minimizing \hat{H}_n over the simplex is a quadratic program for which efficient algorithms are available. The convexity of \hat{H}_n also proves that $\hat{\theta}$ is well defined, although it may not be unique (for example if all $\hat{\mu}_j$ are the same then \hat{H}_n is constant on the simplex).

3.2 Assumptions on the noise

We state here the three different assumptions under which our main result, Theorem 3.1 below, holds. The value of β given below is used in the construction of the estimator $\hat{\theta}$ defined in (3.3). The value of β depends on the assumption on the noise.

The constant $L > 0$ is independent of the noise and its role will be specified in Theorem 3.1.

Assumption 3.1 (Gaussian noise). *Assume that the noise components ξ_1, \dots, ξ_n are normal, independent, zero mean, and ξ_i has variance σ_i^2 . In this case, let*

$$\beta = (12 + 16L + 6L^2) \left(\max_{i=1, \dots, n} \sigma_i^2 \right). \quad (3.8)$$

Assumption 3.2 (Subgaussian noise). *Let $K > 0$ and assume that the noise components ξ_1, \dots, ξ_n are independent, zero mean, $\|\xi_i\|_{\psi_2} \leq K$ and ξ_i has variance σ_i^2 . Here, let*

$$\beta = K^2 \left(c_{w_1} (2 + L) 2L + 2c_h^2 (1 + L)^2 + \frac{1}{2} c_{w_2}^2 \max_{i=1, \dots, n} \frac{\|\xi_i\|_{\psi_2}^2}{\sigma_i^2} (2 + L)^2 \right), \quad (3.9)$$

where c_{w_1}, c_{w_2} and c_h are the absolute constants given in Propositions A.2 and A.3.

Assumption 3.3 (Bernstein condition on ξ_1^2, \dots, ξ_n^2). *Let $K > 0$ and assume that the noise components ξ_1, \dots, ξ_n are independent and satisfy*

$$\forall p \geq 1, \quad \mathbb{E}|\xi_i|^{2p} \leq \frac{1}{2} p! \sigma_i^2 K^{2(p-1)}. \quad (3.10)$$

Here, let

$$\beta = 392 + 1408L + 608L^2. \quad (3.11)$$

Assumption 3.3 is the natural assumption to derive a Bernstein concentration inequality for the sum of random variables $\xi_1^2 + \dots + \xi_n^2$. Although Assumption 3.3 is less common than Assumptions 3.1 and 3.2, its interest resides in the concentration inequality given in Theorem 3.2, which is sharper than the Hanson-Wright inequality. Under this assumption, it is possible to remove the expression $\max_{i=1,\dots,n} \|\xi_i\|_{\psi_2}/\sigma_i$ from the value of β .

3.3 Main result

Theorem 3.1. *Let $L > 0$ be a positive real number and $M \geq 2$. For $j = 1, \dots, M$, consider the estimators $\hat{\mu}_j = A_j \mathbf{y} + b_j$ with $b_j \in \mathbf{R}^n$ and A_j a real matrix of size $n \times n$. Assume that the matrices A_1, \dots, A_M satisfy $\|A_j - A_k\|_2 \leq 2L$ for any j, k .*

Assume one of the Assumptions 3.1, 3.2 or 3.3 on the noise $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n)$ and set the value of β accordingly. Once β is set, let $\hat{\theta}$ be defined in (3.3). Then for all $x > 0$, with probability greater than $1 - 2 \exp(-x)$,

$$\|\hat{\mu}_{\hat{\theta}} - \mathbf{f}\|_2^2 \leq \min_{j=1,\dots,M} \left(\|\hat{\mu}_j - \mathbf{f}\|_2^2 + 2\beta \log \frac{1}{\pi_j} \right) + \beta x. \quad (3.12)$$

The proof of Theorem 3.1 is given in Appendix B.2. We now discuss the assumptions of section 3.2 and compare Theorem 3.1 to previous results.

Subgaussian noise. One of the contribution of the present paper is to provide a sharp oracle inequality such as (3.12) under Subgaussian noise. To our knowledge, (3.12) is the first result on sharp oracle inequality in deviation for model selection type aggregation obtained without assuming that the noise is Gaussian.

The traces of the matrices A_1, \dots, A_M . The sharp oracle inequality in deviation given in [9] presents an additive term proportional to $\sigma^2 \text{Tr}(A_j)$, as in (1.3). An improvement of the present paper is the absence of this additive term which can be large for matrices A_j with large trace. Our analysis shows that the quantities $\sigma^2 \text{Tr}(A_j)$ are not meaningful for the problem of aggregation of affine estimators. So even in the Gaussian noise setting, Theorem 3.1 improves upon the earlier result of [9]. When the uniform prior is used, i.e., $\pi_j = 1/M$ for all $j = 1, \dots, M$, the sharp oracle inequality (3.12) matches the lower bound from [25, Proof of Theorem 5.3 with $S = 1$] showing that 3.12 is optimal in a minimax sense.

Motivation behind Assumption 3.3. Under Assumption 3.2 (Subgaussian noise), our analysis leads to a remainder term that can be large for random variables that have pathologically small variance relatively to their subgaussian norm: β defined in (3.9) is proportional to $\max_{i=1,\dots,n} \|\xi_i\|_{\psi_2}/\sigma_i$. Under Assumption 3.3 which is slightly stronger and prevents the variance from being pathologically small, this issue can be fixed. We will come back to Assumption 3.3 in Section 3.5 below.

The quantities involved in β . The constant β in the oracle inequality is of the order $K^2(1 \vee L^2)$, where K^2 is the supremum of the variances or the supremum of the squared subgaussian norms, and $2L$ upper bounds the operator norms of all $A_j - A_k$ for $j, k = 1, \dots, M$. In most practical cases, L will be smaller than 1 since all admissible estimators of the form $A_j \mathbf{y}$ satisfy $\|A_j\|_2 \leq 1$ [8], thus the fact that β is proportional to $1 \vee L^2$ is not an issue. Interestingly, the operator norm of the matrices A_1, \dots, A_M does not appear in the sharp oracle inequality in expectation given in [11], while it

plays a crucial role here. On the other hand, the factor K^2 may be more problematic, especially for heteroscedastic noise: β is proportional to the largest variance (resp. the largest subgaussian norm) even if most of the noise random variables have small variance (resp. small subgaussian norm).

General matrices A_1, \dots, A_M . We relax all assumptions on the matrices A_1, \dots, A_M , for instance they may be non-symmetric and have negative eigenvalues. Earlier works studied projection matrices [22], assumed some commutativity property of the matrices [11] or their symmetry and positive semi-definiteness [9]. Although it is shown in [8] that all admissible linear estimators are symmetric with non-negative eigenvalues, some linear estimators used in practice are not symmetric. For example, the last example of [11, Section 1.2] (“moving averages”), exhibits linear estimators that need not be symmetric: if two neighbors of the graph i, j have a different number of neighbours, then $a_{ij} \neq a_{ji}$. Our result also shows that the restrictions on the matrices A_1, \dots, A_M present in [22, 11, 9] were not intrinsic to the problem of aggregation of affine estimators.

3.4 Outline of the proof

The following lemma shows that we can derive a sharp oracle inequality for the estimator $\hat{\mu}_{\hat{\theta}}$ by controlling the concentration of terms of the form $\boldsymbol{\xi}^T Q \boldsymbol{\xi}$ and $\boldsymbol{\xi}^T v$, where Q is a $n \times n$ deterministic matrix and v is a deterministic vector in \mathbf{R}^n . We use the following lemma proved in Appendix B.1.

Lemma 3.1. *Let $\hat{\theta}$ be defined in (3.3). Then almost surely,*

$$\|\hat{\mu}_{\hat{\theta}} - \mathbf{f}\|_2^2 \leq \min_{J^*=1, \dots, M} \left(\|\hat{\mu}_{J^*} - \mathbf{f}\|_2^2 + 2\beta \log \frac{1}{\pi_{J^*}} \right) + \max_{j, k=1, \dots, M} \zeta_{j, k}$$

where

$$\begin{aligned} \zeta_{j, k} &= \boldsymbol{\xi}^T Q_{j, k} \boldsymbol{\xi} - \mathbb{E}[\boldsymbol{\xi}^T Q_{j, k} \boldsymbol{\xi}] \\ &\quad + \boldsymbol{\xi}^T v_{j, k} \\ &\quad - \beta \log \frac{1}{\pi_k \pi_j} - \frac{1}{2} \|(A_k - A_j) D_\sigma\|_{\text{HS}}^2 - \frac{1}{2} \|(A_k - A_j) \mathbf{f} + b_k - b_j\|_2^2, \end{aligned} \quad (3.13)$$

and the matrix D_σ , the matrices $Q_{j, k}$ and the vectors $v_{j, k}$ are defined by

$$\begin{aligned} D_\sigma &:= \text{diag}(\sigma_1, \dots, \sigma_n), \\ Q_{j, k} &:= 2(A_k - A_j) - \frac{1}{2}(A_k - A_j)^T (A_k - A_j), \end{aligned} \quad (3.14)$$

$$v_{j, k} := 2 \left(I_{n \times n} - \frac{1}{2}(A_k - A_j)^T \right) ((A_k - A_j) \mathbf{f} + b_k - b_j). \quad (3.15)$$

In Appendix B.2, we prove Theorem 3.1 by applying Lemma 3.1 and controlling the concentration of terms of the form $\boldsymbol{\xi}^T Q_{j, k} \boldsymbol{\xi}$ and $\boldsymbol{\xi}^T v_{j, k}$ under the different Assumptions 3.1, 3.2 and 3.3.

A sketch of the proof of Theorem 3.1 under Assumption 3.1 (Gaussian Noise) goes as follows. The quantity $W_{j, k}^{\text{linear}} := \frac{1}{2} \|(A_k - A_j) \mathbf{f} + b_k - b_j\|_2^2$ in (3.13) is of the order of the variance of $\boldsymbol{\xi}^T v_{j, k}$. Using (A.1) applied to $v = v_{j, k}$, it is shown that for all $t > 0$, with probability greater than $1 - \exp(-t)$,

$$\boldsymbol{\xi}^T v_{j, k} - W_{j, k}^{\text{linear}} \leq \gamma \beta t,$$

where $\gamma \in (0, 1)$ and β is the constant given in (3.8). Similarly, the quantity $W_{j,k}^{\text{quad}} := \frac{1}{2} \|(A_k - A_j)D_\sigma\|_{\text{HS}}^2$ in (3.13) is of the order of the variance of $\boldsymbol{\xi}^T Q_{j,k} \boldsymbol{\xi}$. Using the concentration inequality (A.2) applied to $Q_{j,k}$, we prove that with probability greater than $1 - \exp(-t)$,

$$\boldsymbol{\xi}^T Q_{j,k} \boldsymbol{\xi} - \mathbb{E}[\boldsymbol{\xi}^T Q_{j,k} \boldsymbol{\xi}] - W_{j,k}^{\text{quad}} \leq (1 - \gamma)\beta t.$$

For fixed j and k , these concentration inequalities and the union bound lead to

$$\forall t > 0, \quad \mathbb{P}\left(\zeta_{j,k} + \beta \log \frac{1}{\pi_j \pi_k} > \beta t\right) \leq 2 \exp(-t).$$

Finally, the non-random term $-\beta \log \frac{1}{\pi_k \pi_j}$ is used to perform the union bound on $j, k = 1, \dots, M$, such that for all $x > 0$,

$$\begin{aligned} \mathbb{P}\left(\max_{j,k=1,\dots,M} \zeta_{j,k} > \beta x\right) &\leq \sum_{j,k=1,\dots,M} \mathbb{P}\left(\zeta_{j,k} + \beta \log \frac{1}{\pi_j \pi_k} > \beta(x + \log \frac{1}{\pi_j \pi_k})\right) \\ &\leq \sum_{j,k=1,\dots,M} \pi_j \pi_k 2 \exp(-x) = 2 \exp(-x). \end{aligned}$$

The proof is similar under the two other assumptions 3.2 and 3.3, but different concentration inequalities are used. The proof of Lemma 3.1 can be found in Appendix B.1 and the proof of Theorem 3.1 is given in Appendix B.2.

3.5 Assumption 3.3: examples and concentration inequality

The goal of this section is to present the motivation behind Assumption 3.3 and to present the concentration inequality of Theorem 3.2. This concentration inequality is of independent interest as it provides sharper bounds than the Hanson-Wright inequality.

This assumption is sufficient to remove the quantity $\max_{i=1,\dots,n} \|\xi_i\|_{\psi_2}/\sigma_i$ from the expression (3.9) of β in the sharp oracle inequality of Theorem 3.1. It was the weakest assumption we could find that allowed us to remove the quantity $\max_{i=1,\dots,n} \|\xi_i\|_{\psi_2}/\sigma_i$.

Example 3.1. Centered variables almost surely bounded by K and zero mean Gaussian random variables with variance smaller than K^2 satisfy (3.10).

Example 3.2 (Log-concave random variables). In [27], the authors consider a slightly stronger condition [27, Definition 1.1]. They consider random variables Z satisfying for any integer $p \geq 1$ and some constant K :

$$\mathbb{E}[|Z|^p] \leq p K \mathbb{E}[|Z|^{p-1}], \tag{3.16}$$

and they showed in [27, Section 7] that any distribution that is log-concave satisfies (3.16). Thus, if X^2 is log-concave then our assumption (3.10) holds. See [2, Section 6] for a comprehensive list of the common log-concave distributions.

The next theorem provides a concentration inequality for quadratic forms of independent random variables satisfying the moment assumption (3.10). It is sharper than the Hanson-Wright inequality given in Proposition A.3.

Theorem 3.2. *Assume that the noise random variable $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n)^T$ satisfies Assumption 3.3 for some $K > 0$. Let A be any $n \times n$ real matrix. Then for all $t > 0$,*

$$\mathbb{P}\left(\boldsymbol{\xi}^T A \boldsymbol{\xi} - \mathbb{E}[\boldsymbol{\xi}^T A \boldsymbol{\xi}] > t\right) \leq \exp\left(-\min\left(\frac{t^2}{192K^2 \|AD_\sigma\|_{\text{HS}}^2}, \frac{t}{256K^2 \|A\|_2}\right)\right), \quad (3.17)$$

where $D_\sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$. Furthermore, for any $x > 0$, with probability greater than $1 - \exp(-x)$,

$$\boldsymbol{\xi}^T A \boldsymbol{\xi} - \mathbb{E}[\boldsymbol{\xi}^T A \boldsymbol{\xi}] \leq 256K^2 \|A\|_2 x + 8\sqrt{3}K \|AD_\sigma\|_{\text{HS}} \sqrt{x}. \quad (3.18)$$

The proof of Theorem 3.2 is given in Appendix A.3.2. A key ingredient to prove this concentration result is a decoupling inequality [14, 13]. A simple decoupling inequality for quadratic forms can be found in [31] or [15, Theorem 8.11], and we use this result in order to prove Theorem 3.2.

When t is small, the right hand side of (3.17) becomes

$$\exp\left(-\frac{t^2}{192K^2 \|AD_\sigma\|_{\text{HS}}^2}\right),$$

whereas the right hand side of the Hanson-Wright inequality (A.4) yields

$$\exp\left(-c \frac{t^2}{K^4 \|A\|_{\text{HS}}^2}\right),$$

for some absolute constant $c > 0$. The element of the diagonal matrix D_σ are bounded by K , so Theorem 3.2 gives a sharper bound than the Hanson-Wright inequality in this regime. Under the moment assumption (3.10), we were able to remove the factor $\max_{i=1, \dots, n} (\|\xi_i\|_{\psi_2} / \sigma_i)$ using the concentration inequality from Theorem 3.2.

In particular, the sharp oracle inequality (3.12) with β given in (3.11) holds for all the noise distributions described in Examples 3.1 and 3.2.

4 Sparsity oracle inequality

The goal of this section is to prove a sparsity oracle inequality, when the noise is a subgaussian vector. We are given p deterministic vectors in \mathbf{R}^n that are the columns of a $n \times p$ real matrix \mathbb{X} , and the goal is to find an estimator $\hat{\boldsymbol{\theta}} \in \mathbf{R}^p$ such that the quantity $\|\mathbb{X}\hat{\boldsymbol{\theta}} - \mathbf{f}\|_2^2$ is close to $\|\mathbb{X}\boldsymbol{\theta}^* - \mathbf{f}\|_2^2$ for some sparse $\boldsymbol{\theta}^* \in \mathbf{R}^p$ for which $\mathbb{X}\boldsymbol{\theta}^*$ is a good approximation of the unknown regression vector \mathbf{f} . We will make the following assumption on the noise random vector $\boldsymbol{\xi}$.

Assumption 4.1 (Subgaussian vector). *Let $K > 0$ and assume that the noise random vector $\boldsymbol{\xi}$ satisfy:*

$$\forall \alpha \in \mathbf{R}^n, \quad \mathbb{E} \exp(\alpha^T \boldsymbol{\xi}) \leq \exp\left(\frac{\|\alpha\|_2^2 K^2}{2}\right).$$

Contrary to the previous section, the components of $\boldsymbol{\xi}$ are not assumed to be independent. The same assumption is made in [10]. A direct consequence is the following Hoeffding-type concentration inequality:

$$\mathbb{P}\left(\alpha^T x > K \|\alpha\|_2 \sqrt{2x}\right) \leq \exp(-x). \quad (4.1)$$

Under this assumption, the following concentration inequality was proven in [19].

Proposition 4.1 (One sided concentration [19]). *Let $\boldsymbol{\xi}$ be a random vector in \mathbf{R}^n satisfying Assumption 4.1 for some $K > 0$. Let A be a real $n \times n$ positive semi-definite symmetric matrix. Then for all $x > 0$, with probability greater than $1 - \exp(-x)$,*

$$\boldsymbol{\xi}^T A \boldsymbol{\xi} \leq K^2 (\text{Tr}A + 2 \|A\|_{\text{HS}} \sqrt{x} + 2 \|A\|_2 x). \quad (4.2)$$

This result is remarkable as it holds with the same constants as in the Gaussian case (A.2), under the weak Assumption 4.1. Unlike the previous concentration results given in Appendix A used in Section 3, the above inequality is only one-sided, and it is not known if the above result holds as a two-sided inequality or without the positive semi-definiteness of A . Another difference with the concentration inequalities of Appendix A is that the term $\text{Tr}A$ in (4.2) is an upper bound on the expectation of $\boldsymbol{\xi}^T A \boldsymbol{\xi}$ up to constants. Again, it is not known whether this concentration inequality holds with the constant term $K^2 \text{Tr}A$ replaced by $\mathbb{E}[\boldsymbol{\xi}^T A \boldsymbol{\xi}]$.

The authors of [19] used this result to prove the following inequality for the ordinary least squares estimator $\hat{\mu}_V^{OLS}$ on a d -dimensional linear subspace V of \mathbf{R}^n . The ordinary least squares estimator $\hat{\mu}_V^{OLS}$ is defined as the orthogonal projection of \mathbf{y} on the linear subspace V .

Lemma 4.1 ([19]). *Under Assumption 4.1, with probability greater than $1 - \exp(-x)$:*

$$\begin{aligned} \|\hat{\mu}_V^{OLS} - \mathbf{f}\|_2^2 &\leq \min_{\mu \in V} \|\mu - \mathbf{f}\|_2^2 + K^2(d + 2\sqrt{dx} + 2x), \\ &\leq \min_{\mu \in V} \|\mu - \mathbf{f}\|_2^2 + K^2(2d + 3x). \end{aligned} \quad (4.3)$$

The following corollary extends Proposition 4.1 to general matrices.

Corollary 4.1 (Corollary of Proposition 4.1 for any real matrix A). *Under Assumption 4.1 and for any real matrix A , with probability greater than $1 - \exp(-x)$, the following holds:*

$$\boldsymbol{\xi}^T A \boldsymbol{\xi} \leq K^2 (\|A\|_1 + 2 \|A\|_{\text{HS}} \sqrt{x} + 2 \|A\|_2 x). \quad (4.4)$$

Proof. To see this, let $A_s := \frac{1}{2}(A + A^T)$ and consider $|A_s| := \sqrt{A_s^2}$, the unique square root of the positive semi-definite matrix A_s^2 . By definition of $|A_s|$ and the triangle inequality,

$$\begin{aligned} \boldsymbol{\xi}^T A \boldsymbol{\xi} &= \boldsymbol{\xi}^T A_s \boldsymbol{\xi} \leq \boldsymbol{\xi}^T |A_s| \boldsymbol{\xi}, & \text{Tr}(|A_s|) &= \|A_s\|_1 \leq \|A\|_1, \\ \||A_s|\|_2 &= \|A_s\|_2 \leq \|A\|_2, & \||A_s|\|_{\text{HS}} &= \|A_s\|_{\text{HS}} \leq \|A\|_{\text{HS}}. \end{aligned}$$

Thus applying (4.2) to the matrix $|A_s|$ proves (4.4). \square

Under Assumption 4.1, we obtain the following oracle inequality with proof techniques similar to Theorem 3.1. Define for any $\theta \in \Lambda^M$

$$\begin{aligned} \hat{V}_n(\theta) &= \|\hat{\mu}_\theta\|_2^2 - 2\mathbf{y}^T \hat{\mu}_\theta + K^2 \sum_{j=1, \dots, M} \theta_j (2\|A_j\|_1 + 4\|A_j\|_{\text{HS}}^2) \\ &\quad + \frac{1}{2} \widehat{\text{pen}}(\theta) + \beta \sum_{j=1, \dots, M} \theta_j \log \frac{1}{\pi_j}, \end{aligned}$$

where $\beta > 0$ is a constant, $\widehat{\text{pen}}(\cdot)$ is the penalty (3.2) and the matrix norms are defined in (2.1). We consider the estimator $\hat{\mu}_{\hat{\theta}}$ of \mathbf{f} where

$$\hat{\theta} \in \underset{\theta \in \Lambda^M}{\text{argmin}} \hat{V}_n(\theta). \quad (4.5)$$

The function \hat{V}_n is equal to the the sum of \hat{H}_n (3.1) and some linear function of θ . Thus \hat{V}_n is also convex and minimizing \hat{V}_n over the simplex is a quadratic program.

Proposition 4.2. *Let $K, L > 0$ be real numbers. Assume that the random vector $\boldsymbol{\xi}$ satisfies Assumption 4.1. Assume that the matrices A_1, \dots, A_M satisfy $\|A_j - A_k\|_2 \leq 2L$ for any j, k . Let $\hat{\theta}$ be defined in (4.5) with*

$$\beta = K^2(6 + 8L + 4L^2). \quad (4.6)$$

Then for all $x > 0$, with probability greater than $1 - 2\exp(-x)$,

$$\begin{aligned} \|\hat{\mu}_{\hat{\theta}} - \mathbf{f}\|_2^2 &\leq \min_{j=1, \dots, M} \left(\|\hat{\mu}_j - \mathbf{f}\|_2^2 + 8K^2 \|A_j\|_{\text{HS}}^2 + 4K^2 \|A_j\|_1 \right. \\ &\quad \left. + 2\beta \log \frac{1}{\pi_j} \right) + \beta x. \end{aligned} \quad (4.7)$$

This oracle inequality presents the extra terms proportional to $K^2 \|A_j\|_{\text{HS}}^2$ and $K^2 \|A_j\|_1$ compared to the sharp oracle inequality (3.12). However, this oracle inequality presents some advantages. First, it holds under Assumption 4.1 which is weaker than the noise assumptions of Section 3 since the noise coordinates do not need to be independent. Second, the quantity $\max_{i=1, \dots, n} \|\xi_i\|_{\psi_2} / \sigma_i$ appearing in (3.9) is not present here, which is possible at the cost of the terms proportional to $\|A_j\|_{\text{HS}}^2$ and $\|A_j\|_1$. Finally, one does not need to know the variance of the noise in order to compute the proposed estimator; its construction only relies on K which is an *upper bound* on the subgaussian norm of the random vector $\boldsymbol{\xi}$.

Remark 4.1 (Estimation of an upper bound on the variance). It is easier to construct an estimator that upper bounds the variance than to construct an estimator of the variance itself. For example, for Gaussian noise with variance σ^2 , the estimator $\hat{\sigma}^2$ proposed in [16, Equation (5)] is a positively biased estimator of the variance, for any subspace \mathcal{S}^* of \mathbf{R}^n [16, Section 2.2]. In our setting, a reasonable choice for \mathcal{S}^* is the subspace of vectors proportional to $(1, \dots, 1)$.

We now use the oracle inequality (4.7) to perform sparsity pattern aggregation [25, 24, 9, 28]. For each subset $J \subset \{1, \dots, p\}$, let μ_J^{OLS} be the ordinary least squares

estimator on the linear span of the columns of \mathbb{X} whose indices are in J . This estimator satisfies the oracle inequality (4.3) with $d \leq |J|$, where $|J|$ denotes the cardinal of J . We aggregate these 2^p ordinary least squares estimators using the method (4.5) and the prior distribution given by $\pi_J \propto e^{-|J|} \binom{p}{|J|}^{-1}$. As sparsity pattern aggregation is not central in the present paper, we keep this presentation short and refer the reader to [25, 24, 9, 28] for the construction of ordinary least squares estimators and sparsity pattern aggregation for more details.

As the ordinary least squares estimators are projections of the form $\mu_J^{OLS} = A_J \mathbf{y}$ for some projection matrix A_J , we can take $L = 1$ in Proposition 4.2 and the inequalities $\|A_J\|_1 \leq |J|$ and $\|A_J\|_{\text{HS}}^2 \leq |J|$ hold. Define $\hat{\theta}^{SPA}$ such that

$$\mathbb{X} \hat{\theta}^{SPA} = \hat{\mu}_{\hat{\theta}}, \quad (4.8)$$

where $\hat{\theta}$ is the estimator from (4.5) and $\hat{\mu}_{\hat{\theta}}$ is obtained by aggregating the $M = 2^p$ estimators $(\hat{\mu}_J^{OLS})_{J \subset \{1, \dots, p\}}$. Then the following sparsity oracle inequality holds, where $|\beta|_0$ is the number of non-zero coefficients of β .

Theorem 4.1. *Let \mathbb{X} be a deterministic design matrix with p columns and let $\hat{\theta}^{SPA}$ the sparsity pattern aggregate defined in (4.8). Under Assumption 4.1 on the noise $\boldsymbol{\xi}$, with probability greater than $1 - 3 \exp(-x)$,*

$$\left\| \mathbb{X} \hat{\theta}^{SPA} - \mathbf{f} \right\|_2^2 \leq \inf_{\theta \in \mathbf{R}^p} \left[\|\mathbb{X} \theta - \mathbf{f}\|_2^2 + 21K^2 x \right. \\ \left. K^2 \left(18 + 12|\theta|_0 + 72|\theta|_0 \log \left(\frac{ep}{1 \vee |\theta|_0} \right) \right) \right]. \quad (4.9)$$

Theorem 4.1 improves upon the previous results on sparsity pattern aggregation [9, 25, 24, 28] in several aspects.

First, the noise $\boldsymbol{\xi}$ is only assumed to be subgaussian and its components need not be independent, whereas previous results only hold under Gaussianity and independence of the noise components. Theorem 4.1 shows that the optimal bounds previously known for Gaussian noise [9, 25, 24, 28] are of the same form when the noise is only assumed to be subgaussian.

Second, to construct the aggregates in [9, 25, 24, 28] one needs the exact knowledge of the covariance matrix of the noise. In Theorem 4.1, only an upper bound of the subgaussian norm of the noise is needed to construct the estimator. As explained in Remark 4.1, for Gaussian noise a rough upper bound can be estimated from the data.

Third, we do not split the data in order to perform sparsity pattern aggregation, as opposed to the ‘‘sample cloning’’ approach [28, Lemma 2.1]. Sample cloning is possible only for Gaussian noise when the variance is known; it cannot be used here as $\boldsymbol{\xi}$ can be any subgaussian vector.

The estimator of Theorem 4.1 achieves the minimax rate for any intersection of ℓ_0 and ℓ_q balls, where $q \in (0, 2)$. This can be shown by applying the arguments of [9, 28] and bounding the right hand side of (4.9). Indeed, although [9, 28] consider only normal random variables, the argument does not depend on the noise distribution.

The result above holds without any assumption on the design matrix \mathbb{X} , as opposed to the LASSO or the Dantzig estimators which need assumptions on the design matrix \mathbb{X} to achieve sparsity oracle inequalities similar but weaker than (4.9).

The interest of the LASSO and the Dantzig estimators is that they can be computed efficiently for large p . The sparsity pattern aggregate based on exponential weights can also be computed efficiently using MCMC methods [25]. The estimator $\hat{\theta}^{SPA}$ proposed here suffers the same drawback as [5] or the sparsity pattern aggregate performed with Q -aggregation [9]: it is not known whether these estimators can be computed in polynomial time, which makes them useful only for relatively small p .

A Concentration inequalities

A.1 Gaussian concentration

Let X be a zero mean Gaussian random variable with variance σ^2 . A standard bound on the Gaussian tail is

$$\forall x > 0, \quad \mathbb{P}\left(X > \sigma\sqrt{2x}\right) \leq \exp(-x).$$

Let $v \in \mathbf{R}^n$ and let ξ_1, \dots, ξ_n be zero mean independent Gaussian random variables with $\mathbb{E}[\xi_i^2] = \sigma_i^2$ for all i . Then $v^T \boldsymbol{\xi}$ is Gaussian with variance $\|D_\sigma v\|_2^2$ and thus

$$\forall x > 0, \quad \mathbb{P}\left(v^T \boldsymbol{\xi} > \|D_\sigma v\|_2 \sqrt{2x}\right) \leq \exp(-x). \quad (\text{A.1})$$

Proposition A.1 (Gaussian chaos of order 2). *Let ξ_1, \dots, ξ_n be independent zero mean normal random variables with for all $i = 1, \dots, n$, $\mathbb{E}[\xi_i^2] = \sigma_i^2$. Let A be any $n \times n$ real matrix. Then for any $x > 0$,*

$$\mathbb{P}\left(\boldsymbol{\xi}^T A \boldsymbol{\xi} - \mathbb{E}[\boldsymbol{\xi}^T A \boldsymbol{\xi}] > 2\|D_\sigma A D_\sigma\|_{\text{HS}} \sqrt{x} + 2\|D_\sigma A D_\sigma\|_2 x\right) \leq \exp(-x). \quad (\text{A.2})$$

A proof of this concentration result is given in [6, Example 2.12] for diagonal-free matrices. It can be easily extended to the general case via the following argument.

Proof of Proposition A.1. First, notice that if the result holds for standard normal random variables with variance 1, then by considering the random variables $\xi'_i = \xi_i/\sigma_i$ and the matrix $M = D_\sigma A D_\sigma$, the result also holds when ξ_1, \dots, ξ_n have variances different than 1. Thus in the following we assume without loss of generality that $\sigma_i = 1$ for all $i = 1, \dots, n$.

Second, if the result holds for all symmetric matrices A , then for a non-symmetric matrix A one can consider $B = \frac{A+A^T}{2}$ which is symmetric. Then $\boldsymbol{\xi}^T B \boldsymbol{\xi} = \boldsymbol{\xi}^T A \boldsymbol{\xi}$ and by the triangle inequality,

$$\|B\|_2 \leq \frac{\|A\|_2 + \|A^T\|_2}{2} = \|A\|_2, \quad \|B\|_{\text{HS}} \leq \frac{\|A\|_{\text{HS}} + \|A^T\|_{\text{HS}}}{2} = \|A\|_{\text{HS}}.$$

Thus if the concentration inequality (A.2) holds for the symmetric matrix B , it will also hold for the non-symmetric matrix A . Without loss of generality, we can consider only symmetric matrices.

Let ξ_1, \dots, ξ_n be standard normal random variables and let A be a symmetric matrix. There exists an invertible square matrix U with $U^T = U^{-1}$ such that $A = U^T \Lambda U$ for some diagonal matrix $\Lambda = \text{diag}(\mu_1, \dots, \mu_n)$. By rotational invariance of the normal distribution, if $(X_1, \dots, X_n)^T = U \boldsymbol{\xi}$ then X_1, \dots, X_n are i.i.d. standard normal random variables. As $\mathbb{E}[\boldsymbol{\xi}^T A \boldsymbol{\xi}] = \text{Tr} A = \sum_{i=1}^n \mu_i$,

$$\boldsymbol{\xi}^T A \boldsymbol{\xi} - \mathbb{E}[\boldsymbol{\xi}^T A \boldsymbol{\xi}] = \sum_{i=1}^n \mu_i (X_i^2 - 1).$$

The rest of the proof can be treated exactly as in the proof of [6, Example 2.12], using the bound

$$\forall \lambda \in (0, 1/2), \quad \log \mathbb{E} \exp(\lambda(\xi_i^2 - 1)) \leq \frac{\lambda^2}{1 - 2\lambda},$$

without assuming that A is diagonal-free. □

A.2 Subgaussian concentration

Again, we present tools to control terms of the form $\boldsymbol{\xi}^T Q \boldsymbol{\xi}$ and $v^T \boldsymbol{\xi}$ that appear in Lemma 3.1. Proposition A.2 below provides a concentration result for the latter.

Proposition A.2 (Hoeffding-type inequality [30, Section 5.2.3]). *There exists an absolute constant $C_H > 0$ such that the following holds. Let $n \geq 1$ and ξ_1, \dots, ξ_n be independent zero mean subgaussian random variables with $\max_{i=1, \dots, n} \|\xi_i\|_{\psi_2} \leq K$ for some real number $K > 0$. Let $v \in \mathbf{R}^n$.*

Then for any $x > 0$, with probability greater than $1 - \exp(-x)$,

$$\boldsymbol{\xi}^T v \leq C_H K \|v\|_2 \sqrt{x} \tag{A.3}$$

where $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n)^T$.

The concentration result for a quadratic form of independent zero mean subgaussian random variables given in Proposition A.3 below is known as the Hanson-Wright inequality. First versions of this inequality can be found in Hanson and Wright [18] and Wright [32], although with a weaker statement than Proposition A.3 below since these results involve $\|(|a_{ij}|)\|_2$ instead of $\|A\|_2$. Recent proofs of this concentration inequality with $\|A\|_2$ instead of $\|(|a_{ij}|)\|_2$ can be found in Rudelson and Vershynin [26] or Barthe and Milman [3, Theorem A.5].

Proposition A.3 (Hanson-Wright inequality [26]). *There exist absolute constants $c_{w_1}, c_{w_2}, c > 0$ such that the following holds. Let $n \geq 1$ and ξ_1, \dots, ξ_n be independent zero mean subgaussian random variables with $\max_{i=1, \dots, n} \|\xi_i\|_{\psi_2} \leq K$ for some real number $K > 0$. Let A be any $n \times n$ real matrix. Then for all $t > 0$,*

$$\mathbb{P} \left(\boldsymbol{\xi}^T A \boldsymbol{\xi} - \mathbb{E}[\boldsymbol{\xi}^T A \boldsymbol{\xi}] > t \right) \leq \exp \left(-c \min \left(\frac{t^2}{K^4 \|A\|_{\text{HS}}^2}, \frac{t}{K^2 \|A\|_2} \right) \right) \tag{A.4}$$

where $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n)^T$. Furthermore, for any $x > 0$, with probability greater than $1 - \exp(-x)$,

$$\boldsymbol{\xi}^T A \boldsymbol{\xi} - \mathbb{E}[\boldsymbol{\xi}^T A \boldsymbol{\xi}] \leq c_{w_1} K^2 \|A\|_2 x + c_{w_2} K^2 \|A\|_{\text{HS}} \sqrt{x}. \tag{A.5}$$

A.3 Concentration under Assumption 3.3

A.3.1 Bounds on moment generating functions

The condition (3.10) leads to the following bounds on the moment generating functions of X and X^2 , which are crucial to prove Theorem 3.2.

Proposition A.4. *Let $K > 0$ and let ξ_i be a random variable satisfying (3.10) with $\sigma_i^2 = \mathbb{E}[\xi_i^2]$. Then for all $s \in \mathbf{R}$:*

$$\mathbb{E} \exp(s\xi_i) \leq \exp(s^2 K^2). \quad (\text{A.6})$$

Furthermore, if $0 \leq 2sK^2 \leq 1$, then

$$\mathbb{E} \exp(s\xi_i^2 - s\sigma_i^2) \leq \exp(s^2 \sigma_i^2 K^2), \quad (\text{A.7})$$

$$\mathbb{E} \exp(s\xi_i^2) \leq \exp\left(\frac{3}{2}s\sigma_i^2\right). \quad (\text{A.8})$$

Inequality (A.6) shows that a random variable X satisfying the moment assumption (3.10) is subgaussian and its ψ_2 norm is bounded by K up to a multiplicative absolute constant. For any vector $v \in \mathbf{R}^n$, given n independent variables ξ_1, \dots, ξ_n satisfying the moment assumption (3.10), the following Hoeffding-type inequality holds:

$$\mathbb{P}(v^T \boldsymbol{\xi} > 2K \|v\|_2 \sqrt{x}) \leq \exp(-x), \quad (\text{A.9})$$

it is a direct application of (A.6) combined with the Chernoff bound.

The proof of Proposition A.4 is based on Taylor expansions and some algebra.

Proof of Proposition A.4. To simplify the notation, let $X = \xi_i$ and $\sigma = \sigma_i$. We first prove (A.7). We apply the assumption on the even moments of X :

$$\mathbb{E} \exp(sX^2) = 1 + s\sigma^2 + \sum_{p \geq 2} \frac{s^p \mathbb{E} X^{2p}}{p!} \leq 1 + s\sigma^2 + \frac{\sigma^2 s}{2} \sum_{k=1}^{\infty} (sK^2)^k = 1 + s\sigma^2 + \frac{\sigma^2 K^2 s^2}{2(1 - sK^2)},$$

and using the inequality $0 < 2sK^2 \leq 1$, we obtain:

$$\mathbb{E} \exp(sX^2) \leq 1 + s\sigma^2 + \sigma^2 s^2 K^2 \leq \exp(s\sigma^2 + s^2 \sigma^2 K^2),$$

which completes the proof of (A.7). Inequality (A.8) is a direct consequence of (A.7) after applying again the inequality $2sK^2 \leq 1$.

We now prove (A.6). Using the Cauchy-Schwarz inequality and the assumption on the moments for $p = 2$, we get $\sigma^4 \leq \mathbb{E}[\xi^4] \leq \sigma^2 K^2$, so $\sigma \leq K$. Let $p \geq 1$. For the even terms of the expansion of $\mathbb{E} \exp(sX)$, we get:

$$\frac{s^{2p} \mathbb{E} X^{2p}}{(2p)!} \leq \frac{1}{2} (sK)^{2p} \frac{p!}{(2p)!} \leq \frac{1}{2} \frac{(sK)^{2p}}{p!},$$

where for the last inequality we used $(p!)^2 \leq (2p)!$. For the odd terms, by using the Jensen inequality for $p \geq 1$:

$$\frac{s^{2p+1} \mathbb{E} X^{2p+1}}{(2p+1)!} \leq \frac{s^{2p+1} (\mathbb{E} X^{2p+2})^{\frac{2p+1}{2p+2}}}{(2p+1)!} \leq |sK|^{2p+1} \frac{\left(\frac{(p+1)!}{2}\right)^{\frac{2p+1}{2p+2}}}{(2p+1)!} \leq \frac{1}{2} |sK|^{2p+1} \frac{(p+1)!}{(2p+1)!}.$$

If $|sK| > 1$, we use the inequality $(p+1)!^2 \leq (2p+1)!$ to obtain

$$\frac{s^{2p+1}\mathbb{E}X^{2p+1}}{(2p+1)!} \leq \frac{|sK|^{2(p+1)}}{2((p+1)!)},$$

and by combining the inequality for the even and the odd terms:

$$\begin{aligned} \mathbb{E} \exp(sX) &= 1 + \sum_{p \geq 1} \frac{s^{2p}\mathbb{E}X^{2p}}{(2p)!} + \frac{s^{2p+1}\mathbb{E}X^{2p+1}}{(2p+1)!} \leq 1 + \frac{1}{2} \sum_{p \geq 1} \frac{(sK)^{2p}}{p!} + \frac{|sK|^{2(p+1)}}{(p+1)!} \\ &\leq 1 + \sum_{p \geq 1} \frac{(sK)^{2p}}{p!} = \exp(s^2K^2). \end{aligned}$$

If $|sK| \leq 1$, we use the inequality $(p+1)!p! \leq (2p+1)!$ to obtain

$$\frac{s^{2p+1}\mathbb{E}X^{2p+1}}{(2p+1)!} \leq \frac{(sK)^{2p}}{2(p)!},$$

and by combining the inequality for the even and the odd terms:

$$\begin{aligned} \mathbb{E} \exp(sX) &= 1 + \sum_{p \geq 1} \frac{s^{2p}\mathbb{E}X^{2p}}{(2p)!} + \frac{s^{2p+1}\mathbb{E}X^{2p+1}}{(2p+1)!} \leq 1 + \frac{1}{2} \sum_{p \geq 1} \frac{(sK)^{2p}}{p!} + \frac{(sK)^{2p}}{p!} \\ &= 1 + \sum_{p \geq 1} \frac{(sK)^{2p}}{p!} = \exp(s^2K^2). \end{aligned}$$

□

A.3.2 A concentration inequality for quadratic forms

The goal of this section is to prove Theorem 3.2. We start with preliminary calculations that will be useful in the proof. Let A be any $n \times n$ real matrix. Let $\lambda > 0$ satisfy

$$128\|A\|_2 K^2 \lambda \leq 1, \tag{A.10}$$

and define

$$\eta = 32K^2\lambda^2. \tag{A.11}$$

The inequality (A.10) can be rewritten in terms of η :

$$512K^2\|A\|_2^2\eta \leq 1. \tag{A.12}$$

Let A_0 be the matrix A with the diagonal entries set to 0. Then, using the triangle inequality with $A_0 = A - \text{diag}(a_{11}, \dots, a_{nn})$ and $|a_{ii}| \leq \|A\|_2$ for all $i = 1, \dots, n$, we obtain

$$\|A_0\|_2 \leq 2\|A\|_2. \tag{A.13}$$

Let $B = A_0^T A_0 = (b_{ij})_{i,j=1,\dots,n}$ and let B_0 be the matrix B with the diagonal entries set to 0. Then

$$\forall i = 1, \dots, n, \quad 0 \leq b_{ii} = \sum_{j \neq i} a_{ji}^2 \leq \|A\|_2^2. \tag{A.14}$$

By using the decomposition $B_0 = B - \text{diag}(b_{11}, \dots, b_{nn})$ and the inequality $\|v + v'\|_2^2 \leq 2\|v\|_2^2 + 2\|v'\|_2^2$, (A.14) and (A.13), we have:

$$\begin{aligned} \|B_0 \boldsymbol{\xi}\|_2^2 &\leq 2\|B \boldsymbol{\xi}\|_2^2 + 2 \sum_{i=1}^n b_{ii}^2 \xi_i^2, \\ &\leq 2\|A_0\|_2^2 \|A_0 \boldsymbol{\xi}\|_2^2 + 2\|A\|_2^2 \sum_{i=1}^n b_{ii} \xi_i^2, \\ &\leq 8\|A\|_2^2 \|A_0 \boldsymbol{\xi}\|_2^2 + 2\|A\|_2^2 \sum_{i=1}^n b_{ii} \xi_i^2. \end{aligned}$$

Combining the previous display with (A.12), we obtain for any $K > 0$:

$$16K^2 \eta^2 \|B_0 \boldsymbol{\xi}\|_2^2 \leq (512K^2 \|A\|_2^2 \eta) \left(\frac{\eta}{4} \|A_0 \boldsymbol{\xi}\|_2^2 + \frac{\eta}{16} \sum_{i=1}^n b_{ii} \xi_i^2 \right) \leq \frac{\eta}{4} \|A_0 \boldsymbol{\xi}\|_2^2 + \frac{\eta}{16} \sum_{i=1}^n b_{ii} \xi_i^2. \quad (\text{A.15})$$

Proof of Theorem 3.2. Throughout the proof, let $\lambda > 0$ satisfy (A.10). The value of λ will be specified later.

First we treat the diagonal terms by bounding the moment generating function of

$$S_{\text{diag}} := \sum_{i=1}^n a_{ii} \xi_i^2 - \sum_{i=1}^n a_{ii} \sigma_i^2.$$

Using the independence of ξ_1, \dots, ξ_n and (A.7) with $s = a_{ii} \lambda$ with each $i = 1, \dots, n$:

$$\mathbb{E} \exp(\lambda S_{\text{diag}}) \leq \exp \left(\lambda^2 \sum_{i=1}^n a_{ii}^2 \sigma_i^2 K^2 \right), \quad (\text{A.16})$$

provided that for all $i = 1, \dots, n$, $2|a_{ii}| \lambda K^2 \leq 1$ which is satisfied as (A.10) holds and $|a_{ii}| \leq \|A\|_2$.

Now we bound the moment generating function of the off-diagonal terms. Let

$$S_{\text{off-diag}} := \sum_{i,j=1, \dots, n: i \neq j} a_{ij} \xi_i \xi_j.$$

Let the random vector $\boldsymbol{\xi}' = (\xi'_1, \dots, \xi'_n)^T$ be independent of $\boldsymbol{\xi}$ with the same distribution as $\boldsymbol{\xi}$. We apply the decoupling inequality [31] (see also [15, Theorem 8.11]) to the convex function $s \rightarrow \exp(\lambda s)$:

$$\mathbb{E} \exp(\lambda S_{\text{off-diag}}) \leq \mathbb{E} \exp \left(4\lambda \sum_{i,j=1, \dots, n: i \neq j} a_{ij} \xi'_i \xi_j \right).$$

Conditionally on ξ_1, \dots, ξ_n , for each $i = 1, \dots, n$, we use the independence of ξ'_1, \dots, ξ'_n and (A.6) applied to ξ'_i with $s = 4 \sum_{j=1, \dots, n: i \neq j} a_{ij} \xi_j$:

$$\begin{aligned} \mathbb{E} \exp \left(4\lambda \sum_{i \neq j} a_{ij} \xi'_i \xi_j \right) &\leq \mathbb{E} \exp \left(16K^2 \lambda^2 \sum_{i=1, \dots, n} \left(\sum_{j=1, \dots, n: i \neq j} a_{ij} \xi_j \right)^2 \right), \\ &= \mathbb{E} \exp \left(16K^2 \lambda^2 \|A_0 \boldsymbol{\xi}\|_2^2 \right) = \mathbb{E} \exp \left(\frac{\eta}{2} \|A_0 \boldsymbol{\xi}\|_2^2 \right), \end{aligned}$$

where η is defined in (A.11) and A_0 is the matrix A with the diagonal entries set to 0. Let $B = A_0^T A_0 = (b_{ij})_{i, j=1, \dots, n}$. Then $\|A_0 \boldsymbol{\xi}\|_2^2 = \sum_{i=1}^n b_{ii} \xi_i^2 + \sum_{i \neq j} b_{ij} \xi_i \xi_j$.

We use the Cauchy-Schwarz inequality to separate the diagonal terms from the off-diagonal ones:

$$\left(\mathbb{E} \exp \left(\frac{\eta}{2} \|A_0 \boldsymbol{\xi}\|_2^2 \right) \right)^2 \leq \mathbb{E} \exp \left(\eta \sum_{i=1}^n b_{ii} \xi_i^2 \right) \mathbb{E} \exp \left(\eta \sum_{i \neq j} b_{ij} \xi_i \xi_j \right). \quad (\text{A.17})$$

For the off-diagonal terms of (A.17), using the decoupling inequality [31] (see also [15, Theorem 8.11]) we have:

$$\mathbb{E} \exp \left(\eta \sum_{i \neq j} b_{ij} \xi_i \xi_j \right) \leq \mathbb{E} \exp \left(4\eta \sum_{i \neq j} b_{ij} \xi'_i \xi_j \right).$$

Again, conditionally on ξ_1, \dots, ξ_n , for each $j = 1, \dots, n$, we use (A.6) applied to ξ'_i and the independence of ξ'_1, \dots, ξ'_n :

$$\begin{aligned} \mathbb{E} \exp \left(4\eta \sum_{i \neq j} b_{ij} \xi'_i \xi_j \right) &\leq \mathbb{E} \exp \left(16K^2 \eta^2 \sum_{i=1}^n \left(\sum_{j=1, \dots, n: i \neq j} b_{ij} \xi_j \right)^2 \right), \\ &= \mathbb{E} \exp \left(16K^2 \eta^2 \|B_0 \boldsymbol{\xi}\|_2^2 \right), \\ &\leq \mathbb{E} \exp \left(\frac{\eta}{4} \|A_0 \boldsymbol{\xi}\|_2^2 + \frac{\eta}{16} \sum_{i=1}^n b_{ii} \xi_i^2 \right), \end{aligned}$$

where we used the preliminary calculation (A.15) for the last display. Finally, the Cauchy-Schwarz inequality yields

$$\mathbb{E} \exp \left(4\eta \sum_{i \neq j} b_{ij} \xi_i \xi'_j \right) \leq \sqrt{\mathbb{E} \exp \left(\frac{\eta}{2} \|A_0 \boldsymbol{\xi}\|_2^2 \right)} \sqrt{\mathbb{E} \exp \left(\frac{\eta}{8} \sum_{i=1}^n b_{ii} \xi_i^2 \right)}.$$

We plug this upper bound back into (A.17). After rearranging, we find

$$\left(\mathbb{E} \exp \left(\frac{\eta}{2} \|A_0 \boldsymbol{\xi}\|_2^2 \right) \right)^{3/2} \leq \mathbb{E} \exp \left(\eta \sum_{i=1}^n b_{ii} \xi_i^2 \right) \sqrt{\mathbb{E} \exp \left(\frac{\eta}{8} \sum_{i=1}^n b_{ii} \xi_i^2 \right)}.$$

As $b_{ii} \geq 0$, this implies:

$$\mathbb{E} \exp\left(\frac{\eta}{2} \|A_0 \xi\|_2^2\right) \leq \mathbb{E} \exp\left(\eta \sum_{i=1}^n b_{ii} \xi_i^2\right).$$

For each $i = 1, \dots, n$, we apply (A.8) to the variable ξ_i with $s = b_{ii}\eta \geq 0$. Using the independence of ξ_1^2, \dots, ξ_n^2 , we obtain:

$$\mathbb{E} \exp\left(\eta \sum_{i=1}^n b_{ii} \xi_i^2\right) = \prod_{i=1}^n \mathbb{E} \exp(\eta b_{ii} \xi_i^2) \leq \exp\left(\frac{3}{2}\eta \sum_{i=1}^n b_{ii} \sigma_i^2\right) = \exp\left(\frac{3}{2}\eta \|A_0 D_\sigma\|_{\text{HS}}^2\right).$$

provided that for all $i = 1, \dots, n$, $2K^2 b_{ii} \eta \leq 1$ which is satisfied thanks to (A.10) and (A.14).

We remove η from the above displays using its definition (A.11):

$$\mathbb{E} \exp(\lambda S_{\text{off-diag}}) \leq \exp\left(48\lambda^2 K^2 \|A_0 D_\sigma\|_{\text{HS}}^2\right), \quad (\text{A.18})$$

where A_0 is the matrix A with the diagonal entries set to 0.

Now we combine the bound on the moment generating function of S_{diag} and $S_{\text{off-diag}}$, given respectively in (A.16) and (A.18). Using the Chernoff bound and the Cauchy-Schwarz inequality: we have that for all λ satisfying (A.10),

$$\begin{aligned} \mathbb{P}(S_{\text{diag}} + S_{\text{off-diag}} > t) &\leq \exp(-\lambda t) \mathbb{E}[\exp(\lambda S_{\text{diag}}) \exp(\lambda S_{\text{off-diag}})], \\ &\leq \exp(-\lambda t) \sqrt{\mathbb{E}[\exp(2\lambda S_{\text{diag}})]} \sqrt{\mathbb{E}[\exp(2\lambda S_{\text{off-diag}})]}, \\ &\leq \exp\left(-\lambda t + \lambda^2 K^2 \left(\sum_{i=1}^n \sigma_i^2 a_{ii}^2 + 48 \|A_0 D_\sigma\|_{\text{HS}}^2\right)\right), \\ &\leq \exp\left(-\lambda t + 48\lambda^2 K^2 \|AD_\sigma\|_{\text{HS}}^2\right), \end{aligned} \quad (\text{A.19})$$

where for the last display we used the equality

$$\|AD_\sigma\|_{\text{HS}}^2 = \sum_{i,j=1,\dots,n} a_{ij}^2 \sigma_i^2 = \|A_0 D_\sigma\|_{\text{HS}}^2 + \sum_{i=1}^n a_{ii}^2 \sigma_i^2.$$

It now remains to choose the parameter λ . The unconstrained minimum of (A.19) is attained at $\bar{\lambda} = t/(96K^2 \|AD_\sigma\|_{\text{HS}}^2)$. If $\bar{\lambda}$ satisfies the constraint (A.10), then

$$\mathbb{P}(S_{\text{diag}} + S_{\text{off-diag}} > t) \leq \exp\left(\frac{-t^2}{192K^2 \|AD_\sigma\|_{\text{HS}}^2}\right).$$

On the other hand, if $\bar{\lambda}$ does not satisfy (A.10), then the constraint (A.10) is binding and the minimum of (A.19) is attained at $\lambda_b = 1/(128\|A\|_2 K^2) < \bar{\lambda}$. In this case,

$$-t\lambda_b + \lambda_b^2 48K^2 \|AD_\sigma\|_{\text{HS}}^2 \leq -t\lambda_b + \lambda_b \bar{\lambda} 48K^2 \|AD_\sigma\|_{\text{HS}}^2 = -t\lambda_b + \frac{t}{2}\lambda_b = -\frac{t}{256K^2 \|A\|_2}.$$

Combining the two regimes, we obtain

$$\mathbb{P}(S_{\text{diag}} + S_{\text{off-diag}} > t) \leq \exp\left(-\min\left(\frac{t^2}{192K^2 \|AD_\sigma\|_{\text{HS}}^2}, \frac{t}{256K^2 \|A\|_2}\right)\right).$$

The proof of (3.17) is complete.

Now we prove (3.18). The function

$$t \rightarrow x(t) = \min\left(\frac{t^2}{192K^2 \|AD_\sigma\|_{\text{HS}}^2}, \frac{t}{256K^2 \|A\|_2}\right)$$

is increasing and bijective from the set of positive real numbers to itself. Furthermore, for all $t > 0$,

$$t \leq 8\sqrt{3}K \|AD_\sigma\|_{\text{HS}} \sqrt{x(t)} + 256K^2 \|A\|_2 x(t),$$

so the variable change $x = x(t)$ completes the proof of (3.18). \square

B Proofs

B.1 Proof of Lemma 3.1

We start with some preliminary remarks. If $R_n(\theta) := \|\hat{\mu}_\theta\|_2^2 - 2\mathbf{y}^T \hat{\mu}_\theta$, $R_n(\cdot)$ is differentiable and the following identity holds for any $j = 1, \dots, M$ and $\theta \in \Lambda^M$:

$$\nabla R_n(\theta)^T (e_j - \theta) = \|\hat{\mu}_j - \mathbf{f}\|_2^2 - \|\hat{\mu}_\theta - \mathbf{f}\|_2^2 - 2\xi^T (\hat{\mu}_j - \hat{\mu}_\theta) - \|\hat{\mu}_\theta - \hat{\mu}_j\|_2^2. \quad (\text{B.1})$$

The penalty (3.2) satisfies for any $g \in \mathbf{R}^n$ and any $\theta \in \Lambda^M$:

$$\sum_{j=1}^M \theta_j \|\hat{\mu}_j - g\|_2^2 = \|\hat{\mu}_\theta - g\|_2^2 + \widehat{\text{pen}}(\theta). \quad (\text{B.2})$$

This can be shown by using simple properties of the Euclidean norm, or by noting that the equality above is a bias-variance decomposition. The function $\widehat{\text{pen}}(\cdot)$ is differentiable and for any $j = 1, \dots, M$, and $\theta \in \Lambda^M$, one can check that

$$\begin{aligned} \frac{1}{2} \nabla \widehat{\text{pen}}(\theta)^T (e_j - \theta) &= \frac{1}{2} \|\hat{\mu}_\theta - \hat{\mu}_j\|_2^2 - \frac{1}{2} \widehat{\text{pen}}(\theta), \\ &= \|\hat{\mu}_\theta - \hat{\mu}_j\|_2^2 - \frac{1}{2} \sum_{k=1}^M \theta_k \|\hat{\mu}_j - \hat{\mu}_k\|_2^2, \end{aligned} \quad (\text{B.3})$$

where we used (B.2) with $g = \hat{\mu}_j$ for the last equality.

Proof of Lemma 3.1. Let $J^* = 1, \dots, M$ be a deterministic integer. Since $\hat{\theta}$ minimizes \hat{H}_n over the simplex and \hat{H}_n is convex and differentiable, a simple consequence of the KKT conditions [7, 4.2.3, equation (4.21)] yields:

$$\nabla \hat{H}_n(\hat{\theta})^T (e_{J^*} - \hat{\theta}) \geq 0. \quad (\text{B.4})$$

Let $W := \nabla \hat{H}_n(\hat{\theta})^T(e_{J^*} - \hat{\theta})$. By simple algebraic calculations using (B.1) and (B.3), we have

$$\begin{aligned} W &= \|\hat{\mu}_{J^*} - \mathbf{f}\|_2^2 - \|\hat{\mu}_{\hat{\theta}} - \mathbf{f}\|_2^2 - 2\xi^T(\hat{\mu}_{J^*} - \hat{\mu}_{\hat{\theta}}) \\ &\quad + 2\text{Tr}(D_\sigma A_{J^*} D_\sigma) - \sum_{k=1}^M \hat{\theta}_k \text{Tr}(D_\sigma A_k D_\sigma) \\ &\quad - \frac{1}{2} \sum_{k=1}^M \hat{\theta}_k \|\hat{\mu}_k - \hat{\mu}_{J^*}\|_2^2 + \beta \log \frac{1}{\pi_{J^*}} - \beta \sum_{k=1}^M \hat{\theta}_k \log \frac{1}{\pi_k}. \end{aligned}$$

Since for all $j = 1, \dots, M$, $\text{Tr}(D_\sigma A_j D_\sigma) = \mathbb{E}[\xi^T A_j \xi]$, (B.4) can be rewritten

$$\|\hat{\mu}_{\hat{\theta}} - \mathbf{f}\|_2^2 \leq \|\hat{\mu}_{J^*} - \mathbf{f}\|_2^2 + 2\beta \log \frac{1}{\pi_{J^*}} + Z(J^*, \hat{\theta}), \quad (\text{B.5})$$

where for all $J^* = 1, \dots, M$ and $\theta \in \Lambda^M$,

$$\begin{aligned} Z(J^*, \theta) &:= 2\xi^T(\hat{\mu}_\theta - \hat{\mu}_{J^*}) - 2 \sum_{k=1}^M \theta_k \mathbb{E}[\xi^T (A_k - A_{J^*}) \xi] \\ &\quad - \frac{1}{2} \sum_{k=1}^M \theta_k \|\hat{\mu}_{J^*} - \hat{\mu}_k\|_2^2 - \beta \sum_{k=1}^M \theta_k \log \frac{1}{\pi_k} - \beta \log \frac{1}{\pi_{J^*}}. \end{aligned}$$

The quantity $Z(J^*, \theta)$ is affine in its second argument $\theta \in \Lambda^M$ thus it is maximized at a vertex of Λ^M , and the following upper bounds hold:

$$Z(J^*, \hat{\theta}) \leq \max_{\theta \in \Lambda^M} Z(J^*, \theta) = \max_{k=1, \dots, M} Z(J^*, e_k) \leq \max_{j, k=1, \dots, M} Z(j, e_k). \quad (\text{B.6})$$

Let $\zeta_{j,k} := Z(j, e_k)$ for all $j, k = 1, \dots, M$. From (B.5) and (B.6),

$$\|\hat{\mu}_{\hat{\theta}} - \mathbf{f}\|_2^2 \leq \|\hat{\mu}_{J^*} - \mathbf{f}\|_2^2 + 2\beta \log \frac{1}{\pi_{J^*}} + \max_{j, k=1, \dots, M} \zeta_{j,k},$$

where

$$\zeta_{j,k} = 2\xi^T(\hat{\mu}_k - \hat{\mu}_j) - 2\mathbb{E}[\xi^T (A_k - A_j) \xi] - \frac{1}{2} \|\hat{\mu}_k - \hat{\mu}_j\|_2^2 - \beta \log \frac{1}{\pi_k \pi_j}.$$

Let $B_{jk} = A_k - A_j$, so that $\hat{\mu}_k - \hat{\mu}_j = B_{jk}\xi + (B_{jk}\mathbf{f} + b_k - b_j)$. Then

$$\|\hat{\mu}_k - \hat{\mu}_j\|_2^2 = \|B_{jk}\xi\|_2^2 + \|B_{jk}\mathbf{f} + b_k - b_j\|_2^2 + 2\xi^T B_{jk}^T (B_{jk}\mathbf{f} + b_k - b_j). \quad (\text{B.7})$$

After some algebra, we get

$$\begin{aligned} \zeta_{j,k} &= \xi^T Q_{j,k} \xi - \mathbb{E}[\xi^T Q_{j,k} \xi] + \xi^T v_{j,k} \\ &\quad - \beta \log \frac{1}{\pi_k \pi_j} - \frac{1}{2} \|B_{jk} D_\sigma\|_{\text{HS}}^2 - \frac{1}{2} \|B_{jk}\mathbf{f} + b_k - b_j\|_2^2 \end{aligned}$$

where we used the equality $\|B_{jk} D_\sigma\|_{\text{HS}}^2 = \mathbb{E}[\|B_{jk}\xi\|_2^2]$ and where $Q_{j,k}$ and $v_{j,k}$ are defined in (3.14) and (3.15), respectively. \square

B.2 Proof of Theorem 3.1

Let $K, C_{W_1}, C_{W_2}, C_H > 0$ and a diagonal matrix \bar{D} be parameters that are specified below for each assumption. For any $v \in \mathbf{R}^n$ and any real matrix Q , consider the following concentration inequalities: $\forall x > 0$,

$$\mathbb{P}\left(v^T \boldsymbol{\xi} > C_H K \|v\|_2^2\right) \leq \exp(-x), \quad (\text{B.8})$$

$$\mathbb{P}\left(\boldsymbol{\xi}^T Q \boldsymbol{\xi} - \mathbb{E}[\boldsymbol{\xi}^T Q \boldsymbol{\xi}] > C_{W_2} K \|Q \bar{D}\|_{\text{HS}} \sqrt{x} + C_{W_1} K^2 \|Q\|_2 x\right) \leq \exp(-x). \quad (\text{B.9})$$

Let $(\bar{d}_i)_{i=1, \dots, n}$ be the diagonal elements of the matrix \bar{D} and let

$$\beta = K^2 \left(C_{W_1} (2+L) 2L + 2C_H^2 (1+L)^2 + \frac{1}{2} C_{W_2}^2 \max_{i=1, \dots, n} \frac{\bar{d}_i^2}{\sigma_i^2} (2+L)^2 \right). \quad (\text{B.10})$$

The above concentration inequalities are satisfied under the three assumptions on the noise, with different constants:

- Under Assumption 3.1, set $K = \max_{i=1, \dots, n} \sigma_i$, $\bar{D} = D_\sigma$ and $C_H = \sqrt{2}$, $C_{W_1} = 2$, $C_{W_2} = 2$. With this choice of constants, the value of β (B.10) is equal to the value (3.8), (B.8) becomes exactly (A.1) and (B.9) is a consequence of (A.2) applied to the matrix Q and the random vector $\boldsymbol{\xi}$.
- Under Assumption 3.2, K is given in the assumption, set

$$\bar{D} = \text{diag}(\|\xi_1\|_{\psi_2}, \dots, \|\xi_n\|_{\psi_2}),$$

$C_H = c_h$, $C_{W_1} = c_{w_1}$ and $C_{W_2} = c_{w_2}$ where c_h , c_{w_1} and c_{w_2} are the numerical constants from Propositions A.2 and A.3. With this choice of constants, the value of β (B.10) is equal to the value (3.9), (B.8) becomes exactly (A.3) and (B.9) is a direct consequence of (A.5) applied to the random vector $(\frac{\xi_1}{\|\xi_1\|_{\psi_2}}, \dots, \frac{\xi_n}{\|\xi_n\|_{\psi_2}})$ and the matrix $\bar{D} Q \bar{D}$.

- Under Assumption 3.3, K is given in the assumption, set $\bar{D} = D_\sigma$ and $C_H = 2$, $C_{W_1} = 256$, $C_{W_2} = 8\sqrt{3}$. With this choice of constants, the value of β (B.10) is equal to the value (3.11), (B.8) becomes exactly (A.9) and (B.9) becomes exactly (3.18) applied to the random vector $\boldsymbol{\xi}$ and the matrix Q .

Proof of Theorem 3.1. Let $x > 0$. The concentration inequalities (B.8) and (B.9) always hold, with different constants depending on the assumption on the noise as explained above.

Using Lemma 3.1, it is enough to upper bound $\max_{j,k=1, \dots, M} \zeta_{j,k}$ where $\zeta_{j,k}$ is defined in (3.13). Let $j, k = 1, \dots, M$ be fixed, and let $B_{j,k} = A_k - A_j$. We apply the concentration inequality (B.9) to the matrix $Q_{j,k}$ (3.14) and the concentration inequality (B.8) to the vector $v_{j,k}$ (3.15). With the union bound, on the event where both concentration inequalities hold we get that with probability greater than $1 - 2\exp(-x)$,

$$\begin{aligned} \zeta_{j,k} &\leq C_{W_1} K^2 \|Q_{j,k}\|_2 x + C_{W_2} K \|Q_{j,k} \bar{D}\|_{\text{HS}} \sqrt{x} + C_H K \|v_{j,k}\|_2 \sqrt{x} \\ &\quad - \beta \log \frac{1}{\pi_k \pi_j} - \frac{1}{2} \|B_{j,k} D_\sigma\|_{\text{HS}}^2 - \frac{1}{2} \|B_{j,k} \mathbf{f} + b_k - b_j\|_2^2. \end{aligned} \quad (\text{B.11})$$

Using properties of the operator norm and the Hilbert-Schmidt norm with (3.15), (3.14):

$$\begin{aligned} \|v_{j,k}\|_2 &\leq 2\left(1 + \frac{1}{2}\|B_{j,k}\|_2\right) \|B_{j,k}\mathbf{f} + b_k - b_j\|_2, \\ \|Q_{j,k}\|_2 &\leq \left(2 + \frac{1}{2}\|B_{j,k}\|_2\right) \|B_{j,k}\|_2, \\ \|Q_{j,k}\bar{D}\|_{\text{HS}} &\leq 2\|B_{j,k}\bar{D}\|_{\text{HS}} + \frac{1}{2}\|B_{j,k}^T B_{j,k}\bar{D}\|_{\text{HS}}, \\ &\leq \left(2 + \frac{1}{2}\|B_{j,k}\|_2\right) \|B_{j,k}\bar{D}\|_{\text{HS}} \end{aligned}$$

where we used in the last display that for any square matrices M, C , $\|MC\|_{\text{HS}} \leq \|M\|_2 \|C\|_{\text{HS}}$. We plug these inequalities in (B.11):

$$\begin{aligned} \zeta_{j,k} &\leq \left(2 + \frac{1}{2}\|B_{j,k}\|_2\right) (C_{W_1} K^2 \|B_{j,k}\|_2 x + C_{W_2} K \|B_{j,k}\bar{D}\|_{\text{HS}} \sqrt{x}) \\ &\quad + 2C_H K \left(1 + \frac{1}{2}\|B_{j,k}\|_2\right) \|B_{j,k}\mathbf{f} + b_k - b_j\|_2 \sqrt{x} \\ &\quad - \beta \log \frac{1}{\pi_k \pi_j} - \frac{1}{2} \|B_{j,k} D_\sigma\|_{\text{HS}}^2 - \frac{1}{2} \|B_{j,k}\mathbf{f} + b_k - b_j\|_2^2. \end{aligned}$$

We apply the inequality $st \leq \frac{s^2+t^2}{2}$ twice, first with

$$s = C_{W_2} K \left(2 + \frac{1}{2}\|B_{j,k}\|_2\right) \frac{\|B_{j,k}\bar{D}\|_{\text{HS}}}{\|B_{j,k} D_\sigma\|_{\text{HS}}} \sqrt{x}$$

and $t = \|B_{j,k} D_\sigma\|_{\text{HS}}$, second with $s = 2C_H K \left(1 + \frac{1}{2}\|B_{j,k}\|_2\right) \sqrt{x}$ and $t = \|B_{j,k}\mathbf{f} + b_k - b_j\|_2$. In both cases, the term $\frac{t^2}{2}$ cancels and we obtain

$$\begin{aligned} \zeta_{j,k} &\leq K^2 x \left(C_{W_1} \left(2 + \frac{1}{2}\|B_{j,k}\|_2\right) \|B_{j,k}\|_2 + \frac{1}{2} C_{W_2}^2 \frac{\|B_{j,k}\bar{D}\|_{\text{HS}}^2}{\|B_{j,k} D_\sigma\|_{\text{HS}}^2} \left(2 + \frac{1}{2}\|B_{j,k}\|_2\right)^2 \right) \\ &\quad + 2C_H^2 K^2 \left(1 + \frac{1}{2}\|B_{j,k}\|_2\right)^2 x - \beta \log \frac{1}{\pi_k \pi_j}. \end{aligned}$$

Let $(b_{i,l})_{i,l=1,\dots,n}$ be the elements of the matrix $B_{j,k} = A_k - A_j$, and $(\bar{d}_i)_{i=1,\dots,n}$ be the diagonal elements of the matrix \bar{D} . Since

$$\frac{\|B_{j,k}\bar{D}\|_{\text{HS}}^2}{\|B_{j,k} D_\sigma\|_{\text{HS}}^2} = \frac{\sum_{i,l} \bar{d}_i^2 b_{i,l}^2}{\sum_{i,l} \sigma_i^2 b_{i,l}^2} \leq \max_{i=1,\dots,n} \frac{\bar{d}_i^2}{\sigma_i^2},$$

we obtain $\zeta_{j,k} \leq \beta x - \beta \log \frac{1}{\pi_k \pi_j}$ where β is given in (B.10).

For any $t > 0$, let $x = t + \log \frac{1}{\pi_k \pi_j}$. The inequality $\zeta_{j,k} \leq \beta t$ holds with probability greater than $1 - 2\pi_j \pi_k \exp(-t)$. Using the union bound on $j, k = 1, \dots, M$, we have $\max_{j,k=1,\dots,M} \zeta_{j,k} \leq \beta t$ with probability greater than $1 - \sum_{j,k=1,\dots,M} 2\pi_j \pi_k \exp(-t) = 1 - 2 \exp(-t)$. \square

B.3 Sparsity oracle inequality

Proof of Proposition 4.2. Let $J^* = 1, \dots, M$ be a deterministic integer. Since $\hat{\theta}$ minimizes \hat{V}_n over the simplex and \hat{V}_n is convex and differentiable, a simple consequence of

the KKT conditions [7, 4.2.3, equation (4.21)] yields:

$$\nabla \hat{V}_n(\hat{\theta})^T(e_{J^*} - \hat{\theta}) \geq 0. \quad (\text{B.12})$$

Let $W := \nabla \hat{V}_n(\hat{\theta})^T(e_{J^*} - \hat{\theta})$. Using (B.1), (B.3) and some algebra, we obtain

$$\begin{aligned} W &= \|\hat{\mu}_{J^*} - \mathbf{f}\|_2^2 - \|\hat{\mu}_{\hat{\theta}} - \mathbf{f}\|_2^2 - 2\xi^T(\hat{\mu}_{J^*} - \hat{\mu}_{\hat{\theta}}) \\ &\quad - \frac{1}{2} \sum_{k=1}^M \hat{\theta}_k \|\hat{\mu}_k - \hat{\mu}_{J^*}\|_2^2 + \beta \log \frac{1}{\pi_{J^*}} - \beta \sum_{k=1}^M \hat{\theta}_k \log \frac{1}{\pi_k} \\ &\quad + K^2 \left(4 \|A_{J^*}\|_{\text{HS}}^2 + 2 \|A_{J^*}\|_1 - \sum_{k=1}^M \hat{\theta}_k (4 \|A_k\|_{\text{HS}}^2 + 2 \|A_k\|_1) \right). \end{aligned}$$

Inequality (B.12) can be rewritten as

$$\|\hat{\mu}_{\hat{\theta}} - \mathbf{f}\|_2^2 \leq \|\hat{\mu}_{J^*} - \mathbf{f}\|_2^2 + 8K^2 \|A_{J^*}\|_{\text{HS}}^2 + 4K^2 \|A_{J^*}\|_1 + 2\beta \log \frac{1}{\pi_{J^*}} + z(J^*, \hat{\theta})$$

where

$$\begin{aligned} z(J^*, \hat{\theta}) &:= 2\xi^T(\hat{\mu}_{\hat{\theta}} - \hat{\mu}_{J^*}) - 2K^2 \|A_{J^*}\|_1 - 2K^2 \sum_{k=1}^M \hat{\theta}_k \|A_k\|_1 \\ &\quad - 4K^2 \|A_{J^*}\|_{\text{HS}}^2 - 4K^2 \sum_{k=1}^M \hat{\theta}_k \|A_k\|_{\text{HS}}^2 \\ &\quad - \frac{1}{2} \sum_{k=1}^M \hat{\theta}_k \|\hat{\mu}_k - \hat{\mu}_{J^*}\|_2^2 - \beta \log \frac{1}{\pi_{J^*}} - \beta \sum_{k=1}^M \hat{\theta}_k \log \frac{1}{\pi_k}. \end{aligned}$$

The function $z(J^*, \cdot)$ is affine in its second argument. Thus it is maximized at a vertex of Λ^M , and

$$z(J^*, \hat{\theta}) \leq \max_{\theta \in \Lambda^M} z(J^*, \theta) = \max_{k=1, \dots, M} z(J^*, e_k) \leq \max_{j, k=1, \dots, M} z(j, e_k).$$

As it holds for all deterministic $J^* = 1, \dots, M$, we proved that

$$\begin{aligned} \|\hat{\mu}_{\hat{\theta}} - \mathbf{f}\|_2^2 &\leq \min_{J^*=1, \dots, M} \left(\|\hat{\mu}_{J^*} - \mathbf{f}\|_2^2 + 8K^2 \|A_{J^*}\|_{\text{HS}}^2 + 4K^2 \|A_{J^*}\|_1 \right. \\ &\quad \left. + 2\beta \log \frac{1}{\pi_{J^*}} \right) + \max_{j, k=1, \dots, M} \zeta_{jk}, \end{aligned}$$

where

$$\begin{aligned} \zeta_{jk} &:= z(j, e_k) = 2\xi^T(\hat{\mu}_k - \hat{\mu}_j) - 2K^2(\|A_j\|_1 + \|A_k\|_1) - \beta \log \frac{1}{\pi_j \pi_k} \\ &\quad - \frac{1}{2} \|\hat{\mu}_k - \hat{\mu}_j\|_2^2 - 4K^2(\|A_k\|_{\text{HS}}^2 + \|A_j\|_{\text{HS}}^2). \end{aligned}$$

Let $B_{jk} = A_k - A_j$. Using $\hat{\mu}_k - \hat{\mu}_j = B_{jk}\boldsymbol{\xi} + (B_{jk}\mathbf{f} + b_k - b_j)$ and (B.7), we get

$$\begin{aligned}\zeta_{jk} &= 2\boldsymbol{\xi}^T(A_k - A_j)\boldsymbol{\xi} - 2K^2(\|A_j\|_1 + \|A_k\|_1) - 4K^2(\|A_k\|_{\text{HS}}^2 + \|A_j\|_{\text{HS}}^2) \\ &\quad + \boldsymbol{\xi}^T\alpha_{jk} - \frac{1}{2}\|B_{jk}\mathbf{f} + b_k - b_j\|_2^2 \\ &\quad - \frac{1}{2}\|B_{jk}\boldsymbol{\xi}\|_2^2 - \beta \log \frac{1}{\pi_j\pi_k},\end{aligned}$$

where $\alpha_{jk} := 2(I_{n \times n} - \frac{1}{2}B_{jk}^T)(B_{jk}\mathbf{f} + b_k - b_j)$. The vector α_{jk} satisfies

$$\|\alpha_{jk}\|_2 \leq 2(1 + \frac{1}{2}\|B_{jk}\|_2)\|B_{jk}\mathbf{f} + b_k - b_j\|_2.$$

We have $-\|B_{jk}\boldsymbol{\xi}\|_2^2 \leq 0$ almost surely and by the triangle inequality:

$$\begin{aligned}-2K^2(\|A_j\|_1 + \|A_k\|_1) &\leq -2K^2\|A_k - A_j\|_1, \\ -4K^2(\|A_j\|_{\text{HS}}^2 + \|A_k\|_{\text{HS}}^2) &\leq -2K^2\|A_k - A_j\|_{\text{HS}}^2.\end{aligned}$$

Let $x > 0$. We now apply the concentration inequality (4.4) to the matrix $2(A_k - A_j)$ and the Hoeffding-type inequality (4.1) to the vector α_{jk} . Using the union bound, the following holds with probability greater than $1 - 2\exp(-x)$:

$$\begin{aligned}\zeta_{jk} &\leq 4K^2\|A_j - A_k\|_2x + 4K^2\|A_k - A_j\|_{\text{HS}}\sqrt{x} - 2K^2\|A_k - A_j\|_{\text{HS}}^2 \\ &\quad + 2K(1 + \frac{1}{2}\|B_{jk}\|_2)\|B_{jk}\mathbf{f} + b_k - b_j\|_2\sqrt{2x} - \frac{1}{2}\|B_{jk}\mathbf{f} + b_k - b_j\|_2^2 \\ &\quad - \beta \log \frac{1}{\pi_j\pi_k}.\end{aligned}$$

Finally, using the inequality $st \leq \frac{s^2+t^2}{2}$ we obtain

$$4K^2\|A_k - A_j\|_{\text{HS}}\sqrt{x} - 2K^2\|A_k - A_j\|_{\text{HS}}^2 \leq 2K^2x.$$

We apply this inequality again with $t = \|B_{jk}\mathbf{f} + b_k - b_j\|_2$ and $s = 2K(1 + \frac{1}{2}\|B_{jk}\|_2)\sqrt{2x}$:

$$\begin{aligned}2K(1 + \frac{1}{2}\|B_{jk}\|_2)\|B_{jk}\mathbf{f} + b_k - b_j\|_2\sqrt{2x} - \frac{1}{2}\|B_{jk}\mathbf{f} + b_k - b_j\|_2^2 \\ = st - \frac{t^2}{2} \leq \frac{s^2}{2} = 4K^2(1 + \frac{1}{2}\|B_{jk}\|_2)^2x.\end{aligned}$$

The above displays yield the following bound on ζ_{jk} , with probability greater than $1 - 2\exp(-x)$:

$$\zeta_{jk} \leq K^2(6 + 4\|B_{jk}\|_2 + \|B_{jk}\|_2^2)x - \beta \log \frac{1}{\pi_j\pi_k} \leq \beta(x - \log \frac{1}{\pi_j\pi_k}),$$

since $\beta = K^2(6 + 8L + 4L^2)$. We finish the proof using the change of variable $x' = x - \log \frac{1}{\pi_j\pi_k}$ and the union bound on $j, k = 1, \dots, M$, as in the proof of Appendix B.2. \square

We follow exactly the strategy given in [9, 24] to prove Theorem 4.1 by combining the oracle inequalities (4.3) and (4.7).

Proof of Theorem 4.1. Let $\bar{\theta} \in \mathbf{R}^p$ be a minimizer of the right hand side of (4.9) and let $\bar{J} \subset \{1, \dots, M\}$ be the support of $\bar{\theta}$, hence $|\bar{\theta}|_0 = |\bar{J}|$.

We apply the oracle inequality of Proposition 4.2 with $L = 1$ and the oracle inequality (4.3) to the ordinary least squares $\hat{\mu}_{\bar{J}}^{OLS}$. With the union bound, we have with probability greater than $1 - 3 \exp(-x)$:

$$\begin{aligned} \left\| \mathbb{X} \hat{\theta}^{SPA} - \mathbf{f} \right\|_2^2 &\leq \left\| \hat{\mu}_{\bar{J}}^{OLS} - \mathbf{f} \right\|_2^2 + 8K^2 \|A_{\bar{J}}\|_{\text{HS}}^2 + 4K^2 \|A_{\bar{J}}\|_1 + 2\beta \log \frac{1}{\pi_{\bar{J}}} + \beta x, \\ \left\| \hat{\mu}_{\bar{J}}^{OLS} - \mathbf{f} \right\|_2^2 &\leq \left\| \mathbb{X} \bar{\theta} - \mathbf{f} \right\|_2^2 + K^2 (2|\bar{\theta}|_0 + 3x), \end{aligned}$$

where $A_{\bar{J}}$ is the projection matrix such that $\hat{\mu}_{\bar{J}}^{OLS} = A_{\bar{J}} \mathbf{y}$ and β is given in (4.6). By properties of orthogonal projections, $\|A_{\bar{J}}\|_{\text{HS}}^2 \leq |\bar{\theta}|_0$ and $\|A_{\bar{J}}\|_1 \leq |\bar{\theta}|_0$. Combining the two oracle inequalities above with the following bound from [24, Section 5.2.1]:

$$\log \frac{1}{\pi_{\bar{J}}} \leq 2|\bar{\theta}|_0 \log \left(\frac{ep}{|\bar{\theta}|_0} \right) + \frac{1}{2},$$

it yields

$$\left\| \mathbb{X} \hat{\theta}^{SPA} - \mathbf{f} \right\|_2^2 \leq \left\| \mathbb{X} \bar{\theta} - \mathbf{f} \right\|_2^2 + K^2 12|\bar{\theta}|_0 + \beta + 4\beta |\bar{\theta}|_0 \log \left(\frac{ep}{|\bar{\theta}|_0} \right) + K^2 (3 + \beta)x,$$

and replacing β by the expression of the RHS of (4.6) with $L = 1$ finishes the proof. \square

References

- [1] Sylvain Arlot and Francis R. Bach. Data-driven calibration of linear estimators with minimal penalties. In *Advances in Neural Information Processing Systems*, pages 46–54, 2009.
- [2] Mark Bagnoli and Ted Bergstrom. Log-concave probability and its applications. *Economic theory*, 26(2):445–469, 2005.
- [3] Franck Barthe and Emanuel Milman. Transference principles for log-sobolev and spectral-gap with applications to conservative spin systems. *Communications in Mathematical Physics*, 323(2):575–625, 2013. ISSN 0010-3616. doi: 10.1007/s00220-013-1782-2. URL <http://dx.doi.org/10.1007/s00220-013-1782-2>.
- [4] Pierre Bellec. Optimal exponential bounds for aggregation of density estimators. *arXiv preprint arXiv:1405.3907*, Submitted.
- [5] Lucien Birgé and Pascal Massart. Gaussian model selection. *Journal of the European Mathematical Society*, 3(3):203–268, 2001.
- [6] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford University Press, 2013.

- [7] Stephen Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2009.
- [8] Arthur Cohen. All admissible linear estimates of the mean vector. *The Annals of Mathematical Statistics*, pages 458–463, 1966.
- [9] D. Dai, P. Rigollet, Xia L., and Zhang T. Aggregation of affine estimators. *Electron. J. Stat.*, 8:302–327, 2014.
- [10] Dong Dai, Philippe Rigollet, and Tong Zhang. Deviation optimal learning using greedy Q-aggregation. *The Annals of Statistics*, 40(3):1878–1905, 2012.
- [11] Arnak S Dalalyan and Joseph Salmon. Sharp oracle inequalities for aggregation of affine estimators. *The Annals of Statistics*, 40(4):2327–2355, 2012.
- [12] Arnak S Dalalyan and Alexandre B Tsybakov. Aggregation by exponential weighting and sharp oracle inequalities. In *Learning Theory*, volume 4539 of *Lecture Notes in Computer Science*, pages 97–111. Springer Berlin Heidelberg, 2007. ISBN 978-3-540-72925-9. doi: 10.1007/978-3-540-72927-3_9. URL http://dx.doi.org/10.1007/978-3-540-72927-3_9.
- [13] Victor De la Pena and Evarist Giné. *Decoupling: from dependence to independence*. Springer, New York, 1999.
- [14] Victor H. de la Pena and S. J. Montgomery-Smith. Decoupling inequalities for the tail probabilities of multivariate u -statistics. *The Annals of Probability*, 23(2):806–816, 04 1995. doi: 10.1214/aop/1176988291. URL <http://dx.doi.org/10.1214/aop/1176988291>.
- [15] Simon Foucart and Holger Rauhut. *A mathematical introduction to compressive sensing*. Springer, 2013.
- [16] Christophe Giraud. Mixing least-squares estimators when the variance is unknown. *Bernoulli*, 14(4):1089–1107, 2008.
- [17] Christophe Giraud, Sylvie Huet, and Nicolas Verzelen. High-dimensional regression with unknown variance. *Statistical Science*, 27(4):500–518, 11 2012. doi: 10.1214/12-STS398. URL <http://dx.doi.org/10.1214/12-STS398>.
- [18] David Lee Hanson and Farroll Tim Wright. A bound on tail probabilities for quadratic forms in independent random variables. *The Annals of Mathematical Statistics*, 42(3):1079–1083, 1971.
- [19] Daniel Hsu, Sham Kakade, and Tong Zhang. A tail inequality for quadratic forms of subgaussian random vectors. *Electron. Commun. Probab.*, 17:no. 52, 1–6, 2012. ISSN 1083-589X. doi: 10.1214/ECP.v17-2079. URL <http://ecp.ejpecp.org/article/view/2079>.
- [20] I. M. Johnstone. Function estimation and gaussian sequence models. *Unpublished manuscript*, 2(5.3):2, 2002.

- [21] Guillaume Lecué and Philippe Rigollet. Optimal learning with Q-aggregation. *Ann. Statist.*, 42(1):211–224, 2014.
- [22] Gilbert Leung and Andrew R Barron. Information theory and mixing least-squares regressions. *Information Theory, IEEE Transactions on*, 52(8):3396–3410, 2006.
- [23] Arkadi Nemirovski. Topics in non-parametric statistics. In *Lectures on probability theory and statistics (Saint-Flour, 1998)*, volume 1738 of *Lecture Notes in Mathematics*. Springer, Berlin, 2000.
- [24] P. Rigollet and A. Tsybakov. Sparse estimation by exponential weighting. *Statistical Science*, 27(4):558–575, 2012.
- [25] Philippe Rigollet and Alexandre Tsybakov. Exponential screening and optimal rates of sparse estimation. *The Annals of Statistics*, 39(2):731–771, 2011.
- [26] Mark Rudelson and Roman Vershynin. Hanson-wright inequality and sub-gaussian concentration. *Electron. Commun. Probab.*, 18:no. 82, 1–9, 2013. ISSN 1083-589X. doi: 10.1214/ECP.v18-2865. URL <http://ecp.ejpecp.org/article/view/2865>.
- [27] Warren Schudy and Maxim Sviridenko. Concentration and moment inequalities for polynomials of independent random variables. In *Proceedings of the Twenty-Third Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 437–446. SIAM, 2012.
- [28] A.B. Tsybakov. Aggregation and minimax optimality in high-dimensional estimation. In *Proceedings of the International Congress of Mathematicians*, Seoul, 2014. To appear.
- [29] Alexandre B Tsybakov. Optimal rates of aggregation. In *Learning Theory and Kernel Machines*, pages 303–313. Springer, 2003.
- [30] Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- [31] Roman Vershynin. A simple decoupling inequality in probability theory. 2011. URL <http://www-personal.umich.edu/~romanv/papers/decoupling-simple.pdf>.
- [32] Farrol Tim Wright. A bound on tail probabilities for quadratic forms in independent random variables whose distributions are not necessarily symmetric. *The Annals of Probability*, 1(6):1068–1070, 1973.