



**HAL**  
open science

# On the sign recovery given by the thresholded LASSO and thresholded Basis Pursuit

Patrick J C Tardivel, Maa Lgorzata Bogdan

► **To cite this version:**

Patrick J C Tardivel, Maa Lgorzata Bogdan. On the sign recovery given by the thresholded LASSO and thresholded Basis Pursuit. 2019. hal-01956603v2

**HAL Id: hal-01956603**

**<https://hal.science/hal-01956603v2>**

Preprint submitted on 30 Mar 2019 (v2), last revised 31 Aug 2021 (v7)

**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.

# On the sign recovery given by LASSO, thresholded LASSO and thresholded basis pursuit denoising

Patrick J.C. Tardivel<sup>a\*</sup> and Małgorzata Bogdan<sup>a,b</sup>,

<sup>a</sup> Institute of Mathematics, University of Wrocław, Wrocław, Poland

<sup>b</sup> Lund University, Lund, Sweden

## Abstract

In the high-dimensional regression model  $Y = X\beta^0 + \varepsilon$ , we provide new theoretical results on the probability to recover the sign of  $\beta^0$  by the Least Absolute Selection and Shrinkage Operator (LASSO) and by the thresholded LASSO.

It is well known that the irrepresentability is a necessary condition LASSO to recover the sign of  $\beta^0$  with a large probability. In this article we extend this result by providing a tight upper bound for the probability of LASSO sign recovery. This upper bound is smaller than 1/2 when the irrepresentable condition does not hold and thus generalizes Theorem 2 of Wainwright [28]. The bound is attained when non-null components of  $\beta^0$  tend to infinity and its value, which depends on the tuning parameter  $\lambda$ , is the probability that every null components of  $\beta^0$  is correctly estimated at 0. Thus, this bound can be used to select  $\lambda$  for LASSO, so as to control asymptotically at a given level the family wise error rate: the probability that at least one null component of  $\beta^0$  is not estimate by LASSO at 0.

Irrepresentability is a stringent condition to recover the sign of  $\beta^0$  by LASSO, this condition can be relaxed by filtering out LASSO estimates with an appropriately selected threshold. Indeed, it is well known that LASSO estimates are consistent under weaker conditions than the irrepresentability. In this article we provide new theoretical results in the asymptotic setup under which  $X$  is fixed and non-null components of  $\beta^0$  tend to infinity. Apart from LASSO, our results cover also Basis Pursuit DeNoising (BPDN). Compared to the classical asymptotics, where  $X$  is a  $n \times p$  matrix and both  $n$  and  $p$  tend to  $+\infty$ , our approach allows for reduction of the technical burden. In the result our main theorem takes a simple form:

**When non-null components of  $\beta^0$  are sufficiently large, appropriately thresholded LASSO or thresholded BPDN can recover the sign of  $\beta^*$  if and only if  $\beta^0$  is identifiable with respect to the  $l^1$  norm, i.e.**

$$\text{If } X\gamma = X\beta^0 \text{ and } \gamma \neq \beta^0 \text{ then } \|\gamma\|_1 > \|\beta^0\|_1.$$

---

\*Corresponding author: tardivel@math.uni.wroc.pl

We introduce *irrepresentability* and *identifiability* curves which provide the proportion of  $k$  sparse vectors  $\beta^0$  for which the *irrepresentability* and *identifiability* conditions hold. These curves illustrate that the *irrepresentability* is a much stronger condition than *identifiability* especially when the entries in each row of  $X$  are strongly correlated.

Finally, we illustrate how the knockoff methodology [1, 8] allow to select an appropriate threshold and that thresholded BPDN and LASSO can recover the sign of  $\beta^0$  with a larger probability than adaptive LASSO [32].

**Keywords:** Active set estimation, basis pursuit, Identifiability condition, Irrepresentability condition, LASSO, Sign estimation.

## 1 Introduction

Let us consider the high-dimensional linear model

$$Y = X\beta^0 + \varepsilon, \quad (1)$$

where  $X = (X_1 | \dots | X_p)$  is a  $n \times p$  design matrix, with  $n \leq p$ ,  $\varepsilon$  is a random vector in  $\mathbb{R}^n$ , and  $\beta^0 \in \mathbb{R}^p$  is an unknown vector of regression coefficients. The sign vector of  $\beta^0$  is  $S(\beta^0) = (S(\beta_1^0), \dots, S(\beta_p^0)) \in \{-1, 0, 1\}^p$ , where for  $x \in \mathbb{R}$ ,  $S(x) = \mathbf{1}_{x>0} - \mathbf{1}_{x<0}$ . Our main purpose is to recover  $S(\beta^0)$ . This objective is slightly more general than the aim of recovering the active set  $\text{supp}(\beta^0) := \{i \in \{1, \dots, p\} \mid \beta_i^0 \neq 0\}$ . A natural way to estimate  $S(\beta^0)$  is to take sign of a sparse estimator. The LASSO estimator [26] defined hereafter

$$\widehat{\beta}^L := \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \frac{1}{2} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1 \quad (2)$$

is probably the most famous sparse estimator of  $\beta^0$ .

When  $\text{rank}(X) = n$ , an alternative formulation of LASSO is provided by the Basis Pursuit DeNoising (BPDN) estimator [10], defined as follows

$$\widehat{\beta}^{\text{BPDN}} := \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \|\beta\|_1 \text{ subject to } \|Y - X\beta\|_2^2 \leq R. \quad (3)$$

Given a particular vector  $Y \in \mathbb{R}^n$ , there is a one to one correspondance between the tuning parameter  $\lambda > 0$  and the regularization parameter  $R > 0$ , under which LASSO and BPDN yield the same estimation (see *e.g* page 64 of [18] or the chapter 5.3 of [3]). For example, when  $\lambda = \|X'Y\|_\infty$  and when  $R = \|Y\|_2^2$  then both LASSO and BPDN estimators are equal to  $\mathbf{0}$ . However, the relationship between  $\lambda$  and  $R$  depends on  $Y$  and, in broad generality, given a fixed  $\lambda > 0$  for LASSO, we cannot pick a fixed  $R > 0$  for BPDN under which these both estimators equal. Thus, BPDN and LASSO are not equivalent estimators. The Basis Pursuit (BP) estimator,

solution of (3) when  $R = 0$ , is a particular case of BPDN. As discussed e.g. in [12, 16], BP can be thought of as the limit of LASSO when the tuning parameter  $\lambda$  tends to 0.

## 1.1 Sign recovery by LASSO

Properties of the sign estimator  $S(\widehat{\beta}^L(\lambda)) := (S(\widehat{\beta}_1^L(\lambda)), \dots, S(\widehat{\beta}_p^L(\lambda)))$  (or active set estimator  $\text{supp}(\widehat{\beta}^L(\lambda)) := \{i \in \{1, \dots, p\} \mid \widehat{\beta}_i(\lambda) \neq 0\}$ ) have been intensively studied [17, 21, 28, 31, 32]. Specifically, Zhao and Yu [31] and Zou [32] consider the asymptotic setup under which  $n$  tends to  $+\infty$  and  $p$  is fixed and observe that LASSO can recover  $S(\beta^0)$  only if the restrictive "irrepresentable" condition is fulfilled. These results were further extended to the case of the fixed design matrix  $X$ , where the irrepresentable condition is formulated as follows;

**Definition 1 (Irrepresentability condition)** *Let  $X$  be a  $n \times p$  matrix,  $\beta \in \mathbb{R}^p$ ,  $I := \{i \in \{1, \dots, p\} \mid \beta_i \neq 0\}$ , and  $X_I, X_{\bar{I}}$  be the matrices whose columns are respectively  $(X_i)_{i \in I}$  and  $(X_i)_{i \notin I}$ . The vector  $\beta$  satisfies the irrepresentable condition if  $\ker(X_I) = \mathbf{0}$  and  $\|X_{\bar{I}}' X_I (X_I' X_I)^{-1} S(\beta_I)\|_\infty \leq 1$ .*

According to the Theorem 2 of Wainwright [28], the irrepresentability condition is necessary to recover  $S(\beta^0)$  with high probability. Indeed, when  $\ker(X_I) = \mathbf{0}$ ,  $\|X_{\bar{I}}' X_I (X_I' X_I)^{-1} S(\beta_I^0)\|_\infty > 1$  and both  $\varepsilon$  and  $-\varepsilon$  have the same distribution, then for any selection of the tuning parameter  $\lambda > 0$ ,  $\mathbb{P}(S(\widehat{\beta}^L(\lambda)) = S(\beta^0)) \leq 1/2$ . This result holds also in the noiseless case when  $\varepsilon = \mathbf{0}$ , where the probability to recover  $S(\beta^0)$  is equal to zero. Moreover, Bühlmann and van de Geer [5] (page 192-194) showed that, in the noiseless case, when the irrepresentability strictly holds (*i.e.* when  $\|X_{\bar{I}}' X_I (X_I' X_I)^{-1} S(\beta_I^0)\|_\infty < 1$ ) then the non-random set  $\text{supp}(\beta^L(\lambda))$  recovers  $\text{supp}(\beta^0)$  as soon as non-null components of  $\beta^0$  are sufficiently large. The proof provided in [5] can be easily adapted for the sign recovery.

In this article we provide a new theoretical result on the sign recovery by LASSO. Specifically, Theorem 1 in Section 2 provides an upper bound for the probability of the sign recovery of  $\beta^0$  which depends from  $X, S(\beta^0), \lambda$  and on the distribution of  $\varepsilon$ . This upper bound is attained when non-null components of  $\beta^0$  tend to infinity, in which case it is equal to the limiting probability of an event that  $\text{supp}(\widehat{\beta}^L) \subset \text{supp}(\beta^0)$ . Thus, when non-null components of  $\beta^0$  are sufficiently large and the sparsity is known, this bound can be used to select  $\lambda$  in order to control the Family Wise Error Rate (FWER): the probability that at least one null component of  $\beta^0$  is not estimated at 0.

## 1.2 Sign recovery by thresholded LASSO

It is clear that in the noiseless case, the following identifiability condition is necessary and sufficient to recover  $S(\beta^0)$  by the non-random basis pursuit.

**Definition 2 (Identifiability condition)** *The vector  $\beta \in \mathbb{R}^p$  is identifiable with respect to the design matrix*

$X$  and the  $L_1$  norm (or just identifiable with respect to the  $L_1$  norm) if the following implication holds

$$X\gamma = X\beta \text{ and } \gamma \neq \beta \Rightarrow \|\gamma\|_1 > \|\beta\|_1. \quad (4)$$

Under the identifiability assumption,  $\beta^0$  is sparse. Indeed, Lemma 3 in Tardivel et al. [25] shows that  $k = \text{card}\{i \in \{1, \dots, p\} \mid \beta_i^0 \neq 0\} \leq n$ , i.e.  $\beta^0$  has at least  $p - n$  zeros. On the other hand some assumptions on the sparsity of  $\beta^0$  assure that  $\beta^0$  is identifiable with respect to the  $L_1$  norm. For example when  $\|X_1\|_2 = \dots = \|X_p\|_2 = 1$  and the number of nonzero elements of  $\beta^0$  satisfies the following inequality (called mutual coherence condition)

$$k = \text{card}\{i \in \{1, \dots, p\} \mid \beta_i^0 \neq 0\} \leq \frac{1}{2} \left( 1 + \frac{1}{M} \right), \text{ where } M := \max_{i \neq j} |\langle X_i, X_j \rangle|, \quad (5)$$

then  $\beta^0$  is identifiable with respect to the  $L_1$  norm [14, 18, 20]. In the particular case in which the entries of  $X$  are i.i.d  $\mathcal{N}(0, 1)$  and  $n, p$  are both very large, the phase transition curve of Donoho and Tanner [15] provides, with respect to the undersampling ratio  $n/p \in (0, 1)$ , a bound  $\eta \in (0, 1)$  so that  $\beta^0$  having a sparsity  $k$  is identifiable with respect to the  $L_1$  norm if  $k/n < \eta$ .

According to the Theorem 2, reported in Section 3, for any value of the tuning parameter  $\lambda$  or the regularization parameter  $R$ , the identifiability condition is sufficient and necessary so that LASSO or BPDN estimator allow to separate asymptotically negative, null, and positive components of  $\beta^0$ . This means that, when non null components of  $\beta^0$  is sufficiently large, appropriately thresholded LASSO or BPDN can properly identify the sign of  $\beta^0$  if and only if the identifiability condition holds for  $\beta^0$ .

### 1.3 Graphical illustrations of main results

By definition, the irrepresentable condition depends only on  $S(\beta^0)$  and not on how large are the non null components of  $\beta^0$ . Actually, as claimed in the Proposition 2, the identifiability condition also depends only on  $S(\beta^0)$ . Thus, the comparison of these two conditions can be performed by considering sign vectors in  $\{-1, 0, 1\}^p$ . In Figure 1, the identifiability curve (resp. irrepresentability curve) provides the proportion of sign vectors with  $k$  nonzero elements which satisfy the identifiability condition (resp. irrepresentability condition). Figure 1 illustrates that the identifiability curve is highly above the irrepresentability curve. This observation is not surprising since, according to the proposition 1, the irrepresentable condition implies the identifiability condition. Based on these curves, when non null components of  $\beta^0$  are sufficiently large, one can expect that LASSO allows to recover the sign of  $\beta^0$  when the sparsity  $k$  is smaller than 5, while thresholded LASSO allows to recover the sign of  $\beta^0$  when  $k \leq 20$ .

Figure 2 illustrates Theorem 1, which provides an upper bound for LASSO sign recovery. The bound is reached when non null components of  $\beta^0$  tends to  $+\infty$  and, when the irrepresentable condition holds, the tuning parameter  $\lambda$  can be selected in order to fix this upper bound at an arbitrary level. In this figure, the

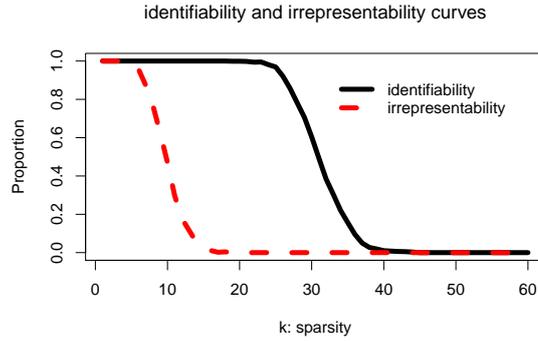


Figure 1: This figure provides the identifiability and irrepresentability curves for the design matrix  $X$  of dimension  $100 \times 300$ , whose entries were independently generated from  $\mathcal{N}(0, 1)$  distribution. The x-axis represents the sparsity  $k$  and the y-axis represents the proportion of sign vectors satisfying the identifiability condition (resp. irrepresentability condition)

design matrix  $X$  is the used in Figure 1 and the noise  $\varepsilon$  is a standard Gaussian vector. According to the irrepresentability curve provided in Figure 1, the irrerepresentable condition holds when  $k = 5$ . The value of  $\lambda$  was selected so that the average value of the bound over 1000 randomly sampled vectors  $\beta^0$  with  $k = 5$  non-zero elements is equal to 0.95. The y axis in Figure 2 represents the probability of recovering  $S(\beta^0)$  calculated based on 1000 randomly sampled vectors  $\beta^0$  having  $k = 5$  non-null components which are equals. Figure 2, shows that the upper bound for LASSO sign recovery is reached when non null components of  $\beta^0$  tend to  $+\infty$  and that the selected  $\lambda$  allows to control the FWER at level 0.05.

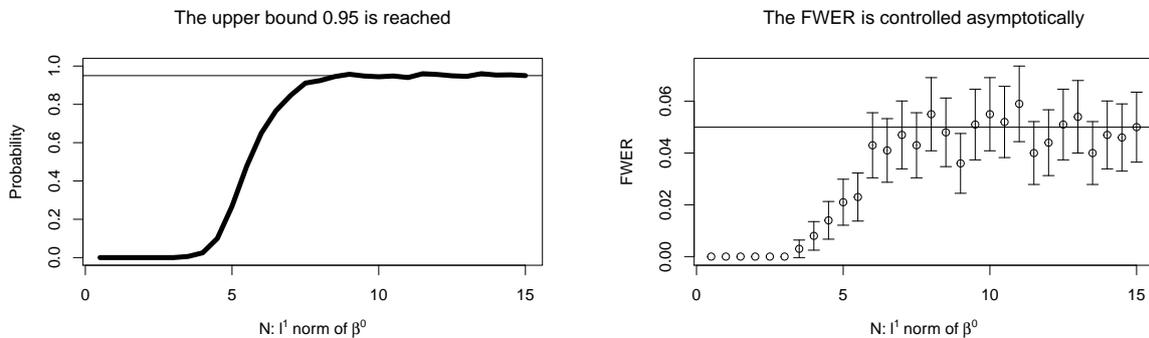


Figure 2: When  $k = 5$ , the left figure provides the probability to recover  $S(\beta^0)$  with LASSO sign estimator and the right figure provides the FWER: the probability that at least one null component of  $\beta^0$  is selected by LASSO estimator. The x-axis represents  $\|\beta^0\|_1$  (in both figures), the y-axis represents the probability of the sign recovery (left figure) and the FWER (right figure). The horizontal lines correspond to  $y = 0.95$  (left figure) and  $y = 0.05$  (right figure).

Figure 3 illustrates Theorem 2 which shows when hen non null components of  $\beta^0$  are large enough, identifiability is a necessary and sufficient condition under which appropriately thresholded BPDN and thresholded LASSO recover  $S(\beta^0)$ . According to the identifiability curve given in the figure 1, when  $k = 20$  the identifiability condition holds. In this figure, the design  $X$  is the same as the one used in Figures 1 and 2 and the y axis

represents the probability of recovering  $S(\beta^0)$  calculated based on 1000 randomly sampled vectors  $\beta^0$  having  $k = 20$  non-null components which are equals. As illustrated in Figure 3, both thresholded BP and thresholded LASSO sign estimators recover  $S(\beta^0)$  when non null components of  $\beta^0$  are large.

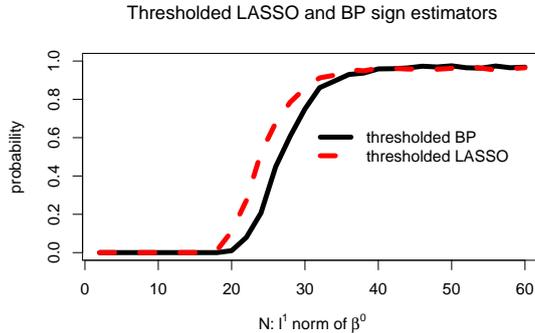


Figure 3: When  $k = 20$ , this figure provides the probability to recover  $S(\beta^*)$  with thresholded LASSO and thresholded BP sign estimators. The x-axis represents  $\|\beta^0\|_1$  and the y-axis represents the sign recovery probability.

Again curves given in figures 2 and 3 are explained with more details in the subsection...

## 1.4 Organization of the article

In section 2, the theorem 1 provides a tight upper bound for LASSO sign estimator to recover  $S(\beta^0)$  and the proposition 1 shows that the irrepresentable condition is stronger than the identifiability condition.

In section 3, the theorem 2 shows that identifiability is a necessary and sufficient condition to recover asymptotically  $S(\beta^0)$  with sign estimators derived from thresholded LASSO and thresholded BPDN.

In section 4, the proposition 2 shows that identifiability condition as irrepresentability condition just depends from  $S(\beta^0)$  and not or on how large are the non-null components of  $\beta^0$ . We introduce the irrepresentability and identifiability curves which provides respectively the proportion of sign vectors satisfying the irrepresentability condition and identifiability condition

The section 5, is devoted to numerical experiments. When  $X$  is a Gaussian matrices with uncorrelated and strongly correlated entries, numerical experiments show that sign estimators derived from the thresholded LASSO and thresholded BPDN are better than both sign estimators derived from LASSO and adaptive LASSO.

## 1.5 Notations and assumptions

In this article we always assume that the design matrix  $X$  is in general condition (see *e.g* [27]; the definition is also reminded in supplementary material). This assumption assures that the minimizer of (2) (resp. minimizer of (3)) is unique and thus that the LASSO estimator (resp. BPDN estimator) is well defined. This assumption is very weak and generically holds. Indeed, when  $X$  is a random matrix such that the entries  $(X_{11}, X_{12}, \dots, X_{np})$

have a density on  $\mathbb{R}^{np}$  then, almost surely,  $X$  is in general position [27].

Hereafter the main notations used in this article:

- Let  $I$  be a subset of  $\{1, \dots, p\}$ , we denote  $\bar{I}$  the complement in  $\{1, \dots, p\}$  of  $I$ , namely  $\bar{I} := \{1, \dots, p\} \setminus I$ .
- The notation  $X_I$  denotes for a matrix whose columns are  $(X_i)_{i \in I}$ .
- Let  $\beta \in \mathbb{R}^p$ , the notation  $\beta_I$  denotes for the vector and  $\text{supp}(\beta)$  denotes for the set  $\{i \in \{1, \dots, p\} \mid \beta_i \neq 0\}$ .
- When  $\beta^0$  is different from  $\mathbf{0}$  we set  $\beta^0 = N\beta^*$  where  $\beta^* = \beta^0 / \|\beta^0\|$  and  $N = \|\beta^0\|$ . In the following  $\beta^*$  is always a fixed unitary vector (for the  $L_1$  norm) and  $N$  is allowed to tends to  $+\infty$  making that non-null components of  $\beta^0$  becomes infinity large.
- LASSO and BPDN estimators depends from  $X, \beta^*, N, \varepsilon$  and from the tuning parameter  $\lambda > 0$  or the regularization parameter  $R \geq 0$ . When it is useful, we add into parentheses these dependencies. The estimator  $\hat{\beta}$  represents indistinctly the LASSO estimator or the BPDN estimator.

## 2 Sign recovery with LASSO sign estimator

The theorem 1 provides an upper bound for the probability to recover  $S(\beta^*)$  with LASSO estimator. When  $\beta^*$  is identifiable with respect to the  $l^1$  norm, this upper bound is reached asymptotically when  $N$  tends to  $+\infty$ .

**Theorem 1** *Let  $I := \text{supp}(\beta^*)$ , let  $X_I, X_{\bar{I}}$  be respectively matrices whose columns are  $(X_i)_{i \in I}$  and  $(X_i)_{i \notin I}$ , let  $\zeta_{X, \lambda, S(\beta^*)}$  be the random vector  $\zeta_{X, \lambda, S(\beta^*)} = X_I' X_I (X_I' X_I)^{-1} S(\beta_I^*) + \frac{1}{\lambda} X_I' (Id - X_I (X_I' X_I)^{-1} X_I') \varepsilon$  and let us assume that  $\ker(X_I) = \mathbf{0}$ .*

**Upper bound:** *The following upper bound, denoted  $\bar{\gamma}$ , for the sign recovery holds.*

$$\mathbb{P} \left( S(\hat{\beta}^{\text{lasso}}(\lambda)) = S(\beta^*) \right) \leq \mathbb{P} \left( \|\zeta_{X, \lambda, S(\beta^*)}\|_{\infty} \leq 1 \right) = \bar{\gamma}.$$

*Now, when  $\beta^*$  is identifiable with respect to the  $l^1$  norm then the following asymptotic results hold.*

**Sharpness of the upper bound:** *Asymptotically, the upper bound is reached.*

$$\begin{aligned} \limsup_{N \rightarrow +\infty} \mathbb{P} \left( S(\hat{\beta}^{\text{lasso}}(\lambda, N)) = S(\beta^*) \right) &\leq \bar{\gamma}, \\ \liminf_{N \rightarrow +\infty} \mathbb{P} \left( S(\hat{\beta}^{\text{lasso}}(\lambda, N)) = S(\beta^*) \right) &\geq \mathbb{P} \left( \|\zeta_{X, \lambda, S(\beta^*)}\|_{\infty} < 1 \right) = \gamma. \end{aligned}$$

**Asymptotic Full power and asymptotic control of the FWER:** *Asymptotically, the power is equal to 1 and the FWER is controlled.*

$$\begin{aligned} \lim_{N \rightarrow +\infty} \mathbb{P} \left( \forall i \in I, S(\widehat{\beta}_i^{\text{lasso}}(\lambda, N)) = S(\beta_i^*) \right) &= 1, \\ \limsup_{N \rightarrow +\infty} \mathbb{P} \left( \exists i \notin I, \widehat{\beta}_i^{\text{lasso}}(\lambda, N) \neq 0 \right) &\leq 1 - \gamma, \\ \liminf_{N \rightarrow +\infty} \mathbb{P} \left( \exists i \notin I, \widehat{\beta}_i^{\text{lasso}}(\lambda, N) \neq 0 \right) &\geq 1 - \bar{\gamma}. \end{aligned}$$

Results given in theorem 1 are quite straightforward when  $X$  is orthogonal (*i.e.* when  $X'X = Id_p$ ). Indeed the upper bound  $\bar{\gamma}$  given in 1) is just the probability that null components of  $\beta^*$  are simultaneously estimated at 0 namely  $\bar{\gamma} = \mathbb{P}(\forall i \notin \text{supp}(\beta^*), \widehat{\beta}_i^{\text{lasso}}(\lambda) = 0)$ . Consequently,  $1 - \bar{\gamma}$  is the Family Wise Error Rate (FWER): the probability that at least one null component of  $\beta^*$  is not estimated at zero.

The famous theorem 2 given in Wainwright [28] is actually a corollary of the upper bound given in 1). Indeed, according to the theorem 1, when  $\varepsilon$  and  $-\varepsilon$  have the same distribution and when  $\|X_I'X_I(X_I'X_I)^{-1}S(\beta_I^*)\|_\infty > 1$  then, whatever  $\lambda > 0$ , the upper bound  $\bar{\gamma}$  is smaller than 1/2.

When  $\gamma = \bar{\gamma}$  (*i.e.* when  $\mathbb{P}(\|\zeta_{X,\lambda,S(\beta^*)}\|_\infty < 1) = \mathbb{P}(\|\zeta_{X,\lambda,S(\beta^*)}\|_\infty \leq 1)$ ) then limits superior and limits inferior given in 2) and 3) imply that

$$\lim_{N \rightarrow +\infty} \mathbb{P} \left( S(\widehat{\beta}^{\text{lasso}}(\lambda, N)) = S(\beta^*) \right) = \bar{\gamma} \text{ and } \lim_{N \rightarrow +\infty} \mathbb{P} \left( \exists i \notin I, \widehat{\beta}_i^{\text{lasso}}(\lambda, N) \neq 0 \right) = 1 - \bar{\gamma}.$$

When the distribution of  $\varepsilon$  and  $S(\beta^*)$  are both known and when the irrepresentable condition strictly holds for  $S(\beta^*)$  then a tuning parameter  $\lambda_{1-\alpha}$  can be chosen such that the upper bound be fixed at  $\bar{\gamma} = 1 - \alpha$ . Since the irrepresentable condition implies the identifiable condition (as proved in the proposition 1), when  $N$  is large, the tuning parameter  $\lambda_{1-\alpha}$  allows to control the FWER at level  $\alpha$ . To our knowledge, the theorem 1, is the first theoretical result providing a guide to select the tuning parameter in order to control a type I error at a specified level.

As claimed above, the proposition 1 shows that the irrepresentable condition on  $\beta^*$  implies the identifiability condition on  $\beta^*$ .

**Proposition 1** *Let  $X$  be a  $n \times p$  matrix with  $n \leq p$  in general position, let  $\beta^* \in \mathbb{R}^p$ , let  $I := \text{supp}(\beta^*)$  and let us assume  $\ker(X_I) = \mathbf{0}$ . If  $\|X_I'X_I(X_I'X_I)^{-1}S(\beta_I^*)\|_\infty \leq 1$ , then the parameter  $\beta^*$  is identifiable with respect to the  $l^1$  norm.*

Let us notice that when the inequality in the irrepresentable condition is strict instead of large, the theorem 1 remains true without assuming that  $X$  is in general position. The proof of the proposition 1 given in this article is the one reported in the PhD manuscript of Tardivel [24].

### 3 Identifiability is a necessary and sufficient condition for sign recovery

When  $\beta^*$  does not satisfy the irrepresentable condition then, even if  $N$  goes to  $+\infty$  and whatever  $\lambda > 0$ , the LASSO sign estimator  $S(\widehat{\beta}^{\text{lasso}}(\lambda))$  fails to recover  $S(\beta^*)$ . However, the irrepresentable condition is not an unsurpassable limitation to recover  $S(\beta^*)$ . Actually the theorem 2 shows that an appropriately thresholded LASSO (resp. thresholded BPDN) recover asymptotically  $S(\beta^*)$  under the identifiability condition on  $\beta^*$  (which is, by the proposition 1, weaker than the irrepresentability condition). To provide a result in a broad generality we do not assume, in the theorem 2, that  $\beta^*$  is identifiable with respect to the  $l^1$  norm. For this theorem, let us introduce the following notations.

- Let  $\tilde{\beta} \in \mathbb{R}^p$  be the solution, in the noiseless case, of the following basis pursuit problem

$$\tilde{\beta} := \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \|\beta\|_1 \text{ subject to } X\beta = X\beta^*.$$

- Let  $\mathcal{B}^- = \{i \in \{1, \dots, p\} \mid \tilde{\beta}_i < 0\}$ ,  $\mathcal{B}^+ = \{i \in \{1, \dots, p\} \mid \tilde{\beta}_i > 0\}$  and  $\mathcal{B} = \mathcal{B}^- \cup \mathcal{B}^+$ .

We remind that,  $\widehat{\beta}$  represents indistinctly the LASSO or BPDN estimator with a fixed tuning parameter  $\lambda > 0$  or with a fixed regularization parameter  $R \geq 0$ .

**Theorem 2** *Let  $X$  be a  $n \times p$  matrix in general position such that  $\operatorname{rank}(X) = n$ .*

**Separation property:** *For any fixed  $\varepsilon \in \mathbb{R}^n$  and sufficiently large  $N > N_0(\varepsilon)$  the following inequality holds*

$$\max_{i \notin \mathcal{B}^-} \left\{ \widehat{\beta}_i(\varepsilon, N) \right\} < \min_{i \notin \mathcal{B}} \left\{ \widehat{\beta}_i(\varepsilon, N) \right\} \leq \max_{i \notin \mathcal{B}} \left\{ \widehat{\beta}_i(\varepsilon, N) \right\} < \max_{i \notin \mathcal{B}^+} \left\{ \widehat{\beta}_i(\varepsilon, N) \right\}.$$

*This inequality means that when  $N$  is large, the estimator  $\widehat{\beta}(\varepsilon, N)$  separates negative components of  $\tilde{\beta}$  (i.e.  $i \in \mathcal{B}^-$ ), null components of  $\tilde{\beta}$  (i.e.  $i \notin \mathcal{B}$ ) and positive components of  $\tilde{\beta}$  (i.e.  $i \in \mathcal{B}^+$ )*

**Sign recovery:** *The equality  $S(\beta^*) = S(\tilde{\beta})$  occurs (thus  $\mathcal{B}^- = \{i \in \{1, \dots, p\} \mid \beta_i^* < 0\}$  and  $\mathcal{B}^+ = \{i \in \{1, \dots, p\} \mid \beta_i^* > 0\}$ ) if and only if  $\beta^*$  is identifiable with respect to the  $l^1$  norm.*

Let us notice that the assumptions on  $X$  are very weak and generically hold when  $n \leq p$ . The assumption  $\operatorname{rank}(X) = n$  assures that, whatever  $R \geq 0$ , the BPDN estimator is well defined. The general position condition assures the uniqueness of both LASSO and BPDN estimators (see e.g. the proposition 1 given in supplementary material for a proof).

Because the almost sure convergence (and thus the convergence for every fixed  $\varepsilon$ ) implies the convergence

in probability then, according to the theorem 2, the following convergence in probability holds

$$\lim_{N \rightarrow +\infty} \mathbb{P} \left( \max_{i \notin \mathcal{B}^-} \{\widehat{\beta}_i(N)\} < \min_{i \notin \mathcal{B}} \{\widehat{\beta}_i(N)\} \leq \max_{i \notin \mathcal{B}} \{\widehat{\beta}_i(N)\} < \max_{i \notin \mathcal{B}^+} \{\widehat{\beta}_i(N)\} \right) = 1.$$

The theorem 2 stress that one cannot recover  $S(\beta^*)$  with a sign estimator derived from LASSO or BPDN when  $\beta^*$  is not identifiable with respect to the  $l^1$  norm since  $S(\beta^*) \neq S(\tilde{\beta})$  (with  $\tilde{\beta}$  as defined in the theorem 2). When  $\beta^*$  is identifiable with respect to the  $l^1$  norm (then  $\tilde{\beta} = \beta^*$ ), the theorem 2 suggest to recover  $S(\beta^*)$  by deriving sign estimators from the thresholded LASSO or thresholded BPDN. Expressions of these thresholded estimators are reported in (6) given below. By the the separation property, one knows that it remains to select a good threshold  $\tau$  to construct a consistent sign estimator (with  $\tau$  depending from  $N$  for the consistency).

The theorem 2 confirms recent results given by Bogdan et al. [4]. Indeed, if  $X$  has i.i.d  $\mathcal{N}(0, 1)$  entries,  $n/p \rightarrow \delta \in (0, 1)$  and if asymptotically the point  $(\text{card}(\text{supp}(\beta^*))/n, n/p)$  is below the asymptotic phase transition curve [13] (*i.e.*  $\beta^*$  is asymptotically identifiable with respect to the  $l^1$  norm) then the thresholded LASSO almost surely recovers  $S(\beta^*)$  (as soon as  $N$  is large enough).

In the following section, we are going to give some properties about identifiability condition. In particular, we show that identifiability (as irrepresentability) just depends from  $S(\beta^*)$  and not on  $N$  or on how large are the non null components of  $\beta^*$ .

## 4 Identifiability and irrepresentability sign applications

By definition the irrepresentable condition just depends from the sign of  $\beta^*$ . Given a particular design matrix  $X$ , the irrepresentability sign application is defined hereafter.

**Irrepresentability sign application:**

$$\Phi_{\text{IC}}^X : s \in \{-1, 0, 1\}^p \mapsto \begin{cases} 1 & \text{if } s = (0, \dots, 0) \\ 1 & \text{if } \ker(X_I) = \mathbf{0} \text{ and } \|X_I' X_I (X_I' X_I)^{-1} s_I\|_\infty \leq 1 \text{ where } I := \text{supp}(s) \\ 0 & \text{otherwise} \end{cases} \quad .$$

Such a sign application provides the limitation of the LASSO sign estimator to recover  $S(\beta^*)$ . Indeed, if  $\phi_{\text{IC}}^X(S(\beta^*)) = 0$  then  $S(\beta^*)$  cannot be recovered with the LASSO sign estimator even if  $N$  is extremely large. The proposition 2 shows that the identifiability condition just depends from  $S(\beta^*)$  and not on  $N$  or on how large are the non null components of  $\beta^*$ .

**Proposition 2** *Let  $X$  be a  $n \times p$  matrix, let  $\beta \in \mathbb{R}^p$  be identifiable with respect to the  $l^1$  norm and let  $\tilde{\beta} \in \mathbb{R}^p$  such that  $S(\beta) = S(\tilde{\beta})$  then  $\tilde{\beta}$  is identifiable with respect to the  $l^1$  norm.*

Given a particular design matrix  $X$ , the identifiability sign application is defined hereafter.

**Identifiability sign application:**

$$\Phi_{\text{Idtf}}^X : s \in \{-1, 0, 1\}^p \mapsto \begin{cases} 0 & \text{if } s \neq \underset{\beta \in \mathbb{R}^p}{\text{argmin}} \|\beta\|_1 \text{ subject to } X\beta = Xs \\ 1 & \text{otherwise} \end{cases}.$$

Such a sign application for the identifiability condition provides the limitation of sign estimators derived from thresholded LASSO and thresholded BPDN to recover  $S(\beta^*)$ . Indeed, if  $\phi_{\text{Idtf}}^X(S(\beta^*)) = 0$  then thresholded LASSO (resp. thresholded BPDN) sign estimator cannot recover  $S(\beta^*)$  even if  $N$  is extremely large.

According to the proposition ... given supplementary material when  $(X_i)_{i \in \text{supp}(\beta^*)}$  is not linearly independent then  $\beta^*$  does not satisfy the identifiability condition. Consequently, when  $\text{card}(\text{supp}(\beta^*)) > n$  then  $\phi_{\text{IC}}^X(S(\beta^*)) = \phi_{\text{Idtf}}^X(S(\beta^*)) = 0$ . Let us provide some basic properties and comments about these sign applications.

1. These two sign applications are even.
2. Due to the proposition 1, whatever  $s \in \{-1, 0, 1\}^p$ ,  $\Phi_{\text{IC}}^X(s) \leq \Phi_{\text{Idtf}}^X(s)$ .
3. The computation of  $\Phi_{\text{IC}}^X$  is a straightforward matricial computation; the computation of  $\Phi_{\text{Idtf}}^X$  is no more difficult and need to solve a basis pursuit problem.

The last remark shows that given a parameter  $\beta^* \in \mathbb{R}^p$ , it is easy to check whether or not  $\beta^*$  is identifiable with respect to the  $l^1$  norm.

Given a sparsity  $k$ , the identifiability (resp. irrepresentability) curve provides the proportion of sign vectors satisfying the identifiability condition (resp. irrepresentability condition). These curves illustrate that the identifiability condition is much weaker than the irrepresentability condition and thus emphasize the theoretical result given in the proposition 1.

#### 4.1 Illustrations of identifiability and irrepresentability curves

The number of sign vectors is very huge ( $3^p$ ), that is why we are not going to provide explicitly  $\Phi_{\text{Idtf}}^X$  and  $\Phi_{\text{IC}}^X$  for each sign vector. Instead, for each sparsity  $k \in \{1, \dots, n\}$ , we are going to compute empirically  $p_{\text{Idtf}}^X(k) := \mathbb{E}_U(\Phi_{\text{Idtf}}^X(U))$  and  $p_{\text{IC}}^X(k) := \mathbb{E}_U(\Phi_{\text{IC}}^X(U))$  where  $U$  is a uniformly distributed on  $\{u \in \{-1, 0, 1\}^p \mid \text{card}(\text{supp}(u)) = k\}$ . The identifiability and irrepresentability curves represent respectively the curves of the functions  $k \in \{1, \dots, n\} \mapsto p_{\text{Idtf}}^X(k)$  and  $k \in \{1, \dots, n\} \mapsto p_{\text{IC}}^X(k)$ . In the numerical experiments given in the figure 4,  $X$  is a Gaussian matrix described hereafter.

**Setting 1:** The matrix  $X$  is a  $n \times p$  matrix with  $n = 100$ ,  $p = 300$  and  $(X_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}$  are i.i.d  $\mathcal{N}(0, 1)$ .

**Setting 2:** The matrix  $X$  is a  $n \times p$  matrix with  $n = 100$ ,  $p = 300$  and the vectors  $(X_{ij})_{1 \leq j \leq p}$  where  $i \in$

$\{1, \dots, n\}$  is a family of i.i.d Gaussian vector  $\mathcal{N}(\mathbf{0}, \Gamma)$ . In this setting  $\Gamma$  is a  $p \times p$  matrix where  $\Gamma_{ii} = 1$  with  $i \in \{1, \dots, p\}$  and  $\Gamma_{ij} = 0.9$  when  $i \neq j$ .

From now on,  $X$  is a particular observation of Gaussian matrix as described in setting 1 and setting 2 (by using the R command `set.seed(123)`).

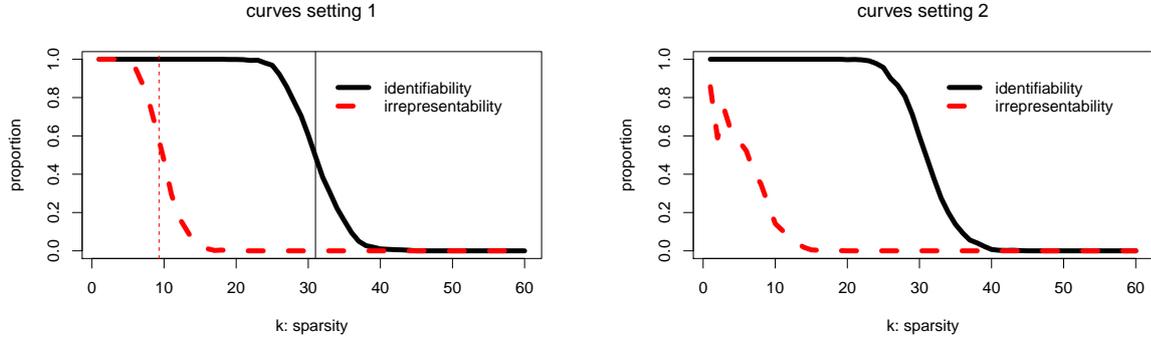


Figure 4: This figure gives the curves of the functions  $k \mapsto p_{\text{Idtf}}^X(k)$  and  $k \mapsto p_{\text{IC}}^X(k)$  when  $X$  is a Gaussian matrix given in the setting 1 (left panel) and setting 2 (right panel). Due to the proposition 1, whatever the sparsity  $k$ ,  $p_{\text{Idtf}}^X(k) \geq p_{\text{IC}}^X(k)$  thus this figure just emphasizes that the identifiability condition is a much weaker assumption than the irrepresentability condition. The vertical lines in the left panel provides, in the setting 1, an asymptotic approximation of the identifiability and irrepresentability curves. Indeed by the theorem 1 in [15] and the theorem 1 in [28], when  $p$  is very large and  $n/p = 1/3$  then the identifiability and irrepresentability conditions hold respectively when  $k \leq 0.31n$  and  $k \leq 0.09n$ . To plot these these curves, for a sparsity  $k$  the quantities  $p_{\text{Idtf}}^X(k)$  and  $p_{\text{IC}}^X(k)$  have been computed by simulating 1000 observations of the random vector  $U$ .

Surprisingly the two identifiability curves given in the setting 1 and 2 are very similar. *A priori*, we expected to recover a curve in the setting 2 much below than the one given in the setting 1. Indeed, classical conditions implying the identifiability of  $\beta^*$  with respect to the  $l^1$  norm are the mutual coherence condition (5) and the restricted isometry property [6, 7]. These conditions are quite weak when the family  $(X_i)_{1 \leq i \leq p}$  is almost orthogonal (as in the setting 1 since  $\mathbb{E}(X'X) = nId_n$ ) but are very strong when  $(X_i)_{1 \leq i \leq p}$  is far from an orthogonal family (as in the setting 2 since  $\mathbb{E}(X'X) = n\Gamma$ ).

The asymptotic phase transition given in Donoho and Tanner [15] provides an approximation of the identifiability curve in the setting 1. Such an approximation is useful when  $n$  and  $p$  are too much large so that the identifiability curve is too much time expensive to obtain. Unfortunately, to our knowledge, there is not such asymptotic phase transition curve for Gaussian matrices with correlated entries as in the setting 2 (see *e.g.* [22] for more details about asymptotic phase transition curve).

One notices that in the setting 2, the irrepresentability curve is not monotonic in the neighbourhood of 0; it is not a numerical problem. Actually when  $k$  is very small, components of  $U$  are all positive or all negative with a quite large probability. Furthermore the figure 5 illustrates that, in the setting 2, when the sign vector  $s$  is positive componentwise (resp. negative componentwise), the irrerepresentable condition becomes a very strong condition. These both remarks, aim at explaining why, in the setting 2, the irrerepresentability curve is not

monotonic. Hereafter, without any loss of generality, we focus on the particular case in which sign vector is positive componentwise. The figure 5 provides the positive irrepresentability and identifiability curves, which are respectively the curves of the functions  $k \mapsto p_{\text{Idtf}+}^X(k) := \mathbb{E}_U(\Phi_{\text{Idtf}}^X(U))$  and  $k \mapsto p_{\text{IC}+}^X(k) := \mathbb{E}_U(\Phi_{\text{IC}}^X(U))$  where  $U$  has uniform distribution over the set  $\{u \in \{0, 1\}^p \mid \text{card}(\text{supp}(u)) = k\}$ .

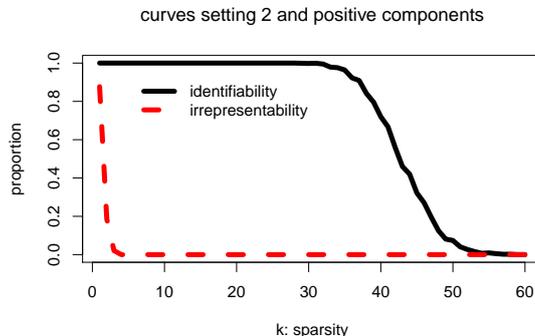


Figure 5: This figure gives the curves of the functions  $k \mapsto p_{\text{Idtf}+}^X(k)$  and  $k \mapsto p_{\text{IC}+}^X(k)$  when  $X$  is a Gaussian matrix given in setting 2. One notices that, with respect to the curves given in the figure 4, the gap between the irrepresentable condition and the identifiability condition becomes larger. When  $k$  is small  $p_{\text{IC}+}^X(k) \approx p_{\text{IC}}^X(k)$  (more precisely,  $p_{\text{IC}+}^X(k) = p_{\text{IC}}^X(k)$  when  $k = 1$ ) and when  $k$  is large enough  $p_{\text{IC}}^X(k)$  weakly depends from the correlation. This remark aim to explain why, in the setting 2, the function  $k \mapsto p_{\text{IC}}^X(k)$  is not monotonic in the neighbourhood of 0. To plot these these curves, for a sparsity  $k$ , the quantities  $p_{\text{Idtf}+}^X(k)$  and  $p_{\text{IC}+}^X(k)$  have been computed by simulating 1000 observations of  $U$ .

Performance of the sign estimators derived from LASSO, thresholded LASSO and thresholded BPDN depends from the tuning parameter  $\lambda$ , regularization parameter  $R$  and threshold  $\tau$ . In the following section, we are going to prescribe values for these parameters.

## 5 Numerical comparisons of sign estimators

Theorem 2 states that the sign estimators provided by thresholded LASSO or thresholded basis pursuit allow to recover  $\text{sign}(\beta^0)$  as long as the identifiability condition is satisfied. Another way to recover  $\text{sign}(\beta^0)$  is to use a sign estimator derived from adaptive LASSO. Indeed, as claimed by theorem 2 of Zou [32], by deriving weights for adaptive LASSO from a consistent estimator of  $\beta^0$  and by selecting properly the tuning parameter  $\lambda$ , one obtains a sign estimator derived from adaptive LASSO which is consistent for  $\text{sign}(\beta^0)$ . Weights for adaptive LASSO can be appropriately derived from LASSO. Indeed, according to lemma 1, under the identifiability assumption LASSO estimator converges to  $\beta^0$ . The purpose of this section is to provide a numerical comparison of sign estimators derived from LASSO, thresholded LASSO, thresholded BP and adaptive LASSO.

As explained hereafter, there are some recommendations on how to select the tuning parameter  $\lambda > 0$  for the LASSO estimator as described in (2) whereas, to our knowledge, there are not clear recommendations on how to select  $R \geq 0$  for BPDN estimator described in (3). That is the reason why, we arbitrary set  $R = 0$  and thus we only consider BP estimator.

### 5.1 Selection of the tuning parameter

As explained in [30, 4], a value of the optimal tuning parameter for the sign recovery by thresholded LASSO is substantially smaller than the optimal value of the tuning parameter for vanilla LASSO. Specifically:

- For LASSO sign estimator, the tuning parameter has to be large enough so that it prevents including false discoveries.
- For thresholded LASSO sign estimator the tuning parameter needs to be selected so as to minimize the mean square error of the estimation of  $\beta^0$ . This tuning parameter does not need to be large, since the threshold will allow to eliminate false discoveries.

#### 5.1.1 Tuning parameter for LASSO sign estimator

When  $\beta^0$  has a known sign so that  $S(\beta^0)$  satisfies the irrepresentable condition, by proposition 1, one may pick a tuning parameter  $\lambda_L$  so that  $\mathbb{P}(S(\hat{\beta}(\lambda_L)) = S(\beta^0))$  is smaller than a given value (say 0.95). Now, according to the irrepresentability curve for our matrix  $X$  with independent columns, the irrepresentability condition is satisfied with probability close to 1 if  $\beta^0$  contains  $k = 5$  nonzero elements. Thus in this setting, we can chose  $\lambda_L$  such that the average value of the upper-bound given in proposition 1 is equal to 0.95. In other words,  $\lambda_L$  is chosen so that  $\mathbb{E}_S(\zeta_{X,\lambda_L,S}) = 0.95$ , where  $S$  is a random sign vector having a uniform distribution over the set  $\{s \in \{-1, 0, 1\}^p \mid \text{card}(\text{supp}(s)) = 5\}$ . The computation of this value gives  $\lambda_L = 8.118$ . Since under the remaining scenarios of our simulation study the irrepresentability condition is typically not satisfied and FWER can not be controlled at the low level, we decided to use the same value  $\lambda_L = 8.118$  for all our simulations.

### 5.1.2 Tuning parameter for thresholded LASSO sign estimator

When  $X$  is the gaussian matrix with independent entries the tuning parameter was selected with the help of the asymptotic theory of Approximate Message Passing (AMP) algorithm for LASSO, provided e.g. in [2, 4, 23]. In the set-up of this theory the design matrix is Gaussian with i.i.d  $\mathcal{N}(0, 1/\sqrt{n})$  and components of  $\beta^0$  are i.i.d random variables having  $\Pi^*$ :  $\Pi = (1 - \gamma)\delta_0 + \gamma\Pi^*$  mixture distribution, where  $\delta_0$  and  $\Pi$  are point mass at 0 distribution and an arbitrary distribution. The number of observation  $n$ , the number of explanatory variables  $p$  becomes infinity large and  $n/p \rightarrow \delta > 0$ . The tuning parameter  $\lambda_{AMP}$ , depending on  $\delta, \gamma, \Pi^*$ , is selected so as to minimize the asymptotic mean square error according to the prescription provided in [2, 30, 4]. As discussed in [30, 4], for any fixed type I error, such a tuning parameter allows to maximize the power. In practice, to compute  $\lambda_{AMP}$ , we replaced the asymptotic parameters of the AMP theory with their finite sample counterparts. Namely,  $\delta = n/p = 100/300$ ,  $\gamma = k/p = k/300$  and  $\Pi^* = \Pi^* = 1/2\delta_t + 1/2\delta_{-t}$ , where  $\delta_t$  is a point mass distribution at  $t$ . Given these parameters, the formula to evaluate  $\lambda_{AMP}$  is provided e.g. in [2, 4, 23]. In case of strongly correlated design we additionally use  $\lambda_s = 0.5\lambda_{AMP}$ .

## 5.2 Selection of the threshold

We aim to construct a sign estimator derived from the thresholded LASSO estimator (resp. thresholded BP estimator) as defined hereafter

$$\forall i \in \{1, \dots, p\}, \widehat{\beta}_i^\tau := \widehat{\beta}_i \mathbf{1}_{\{|\widehat{\beta}_i| > \tau\}} \quad (6)$$

By taking  $\tau_{1-\alpha}$  as the  $1 - \alpha$  quantile of  $\max \left\{ \left| \widehat{\beta}_i \right|, i \notin \text{supp}(\beta^0) \right\}$  (resp.  $\max \left\{ \left| \widehat{\beta}_i \right|, i \notin \text{supp}(\beta^0) \right\}$ ) then the probability to estimate simultaneously every null components of  $\beta^0$  at zero is  $1 - \alpha$ . Consequently, using the threshold  $\tau_{1-\alpha}$  and when non-null components of  $\beta^0$  are very large then thresholded LASSO sign estimator (resp. thresholded BP estimator) recovers  $S(\beta^0)$  with a probability arbitrarily close to  $1 - \alpha$ . Obviously  $\tau_{1-\alpha}$  cannot be obtained by a straightforward computation since  $\beta^0$  is not known.

Given a threshold  $\tau > 0$ , let us set the FWER as follows (the FWER can be seen as the Family Wise Error Rate for multiple testing procedure)

$$\text{FWER} := \mathbb{P} \left( \exists i \notin \text{supp}(\beta^0), \left| \widehat{\beta}_i^\tau \right| \neq 0 \right).$$

In order to provide a threshold larger than  $\tau_{1-\alpha}$  (and thus to control the FWER at level  $\alpha$ ), it could seem appealing to look at the distribution of supremum norm of the LASSO (resp. BP estimator) in the full null model when  $\beta^0 = \mathbf{0}$  [19]. For the BP estimator, Descloux and Sardy [12] suggest the threshold  $\tau_{1-\alpha}^{\text{fn}}$  defined as the  $1 - \alpha$  quantile of  $\max \left\{ \left| \widehat{\beta}_1^{\text{fn}} \right|, \dots, \left| \widehat{\beta}_p^{\text{fn}} \right| \right\}$  where  $\widehat{\beta}^{\text{fn}}$  is the following estimator

$$\widehat{\beta}^{\text{fn}} := \text{argmin} \|\beta\|_1 \text{ subject to } X\beta = \varepsilon.$$

Unfortunately, in the high-dimensional linear model, this intuitive method provides a threshold  $\tau_{1-\alpha}^{\text{fn}}$  which is smaller than  $\tau_{1-\alpha}$  and thus does not assure that  $\text{FWER} \leq \alpha$  (see also Su et al. [23] for additional explanations).

Recently developed knockoff methodology [1, 8], allows to approximate the distribution of  $\widehat{\beta}(\lambda)$  associated to null-components of  $\beta^0$  by creating fake copies of explanatory variables. Consequently, the knockoff methodology is useful to compute a threshold. For this numerical study, we use model free knockoffs proposed in [8] to recover a threshold which heuristically control the FWER at a given level. The approach developed hereafter is available when  $X$  is a Gaussian matrix having a distribution invariant by columns' permutation. In this setting, the size of the knockoff matrix can be as small as possible (see Weinstein et al. [29] for a similar approach). Because adding some fake copies of explanatory variables can change some relevant properties (such as the identifiability condition for  $\beta^0$ ), ideally the knockoff matrix should have just one column. Specifically, at the first step we use model free knockoffs [8] to generate  $30 = p/10$  of fake variables. Then Lasso or BP is run on the matrix supplemented with these additional columns and the maximum of the absolute values of regression coefficients over 30 fake variables is saved. This step is repeated 10 times and the overall maximum of the  $p = 300$  absolute values of regression coefficients over fake variables is calculated. The whole procedure is repeated many times (here 1000) and 0.95 quantile of the obtained maxima is used as the threshold to identify null-components of  $\beta^0$ .

To confirm with the set-up of simulations used to derive the irrepresentability and identifiability curves, in all replicates of our simulation study we used the same fixed design matrix  $X$  described in settings 1 and 2 of the subsection .... In our numerical experiments we randomly sampled the location of the true signals and we randomly generated the error term.

### 5.2.1 LASSO and Adaptive LASSO

In our numerical experiments we selected the following values of the tuning parameters for LASSO and adaptive LASSO:

- For LASSO we selected  $\lambda_L = 8.118$ .
- For the adaptive LASSO the weights are derived using initial estimates  $\widehat{\beta}^L(\lambda_{AMP})$ , where the tuning parameter is selected according to AMP theory, described above. For  $i \in \{1, \dots, p\}$ , weights  $w(\beta_i)$  are defined as  $w(\beta_i) := 1/(\widehat{\beta}_i^L(\lambda_{AMP}) + 10^{-7})$ . Using these weights and the tuning parameter  $\lambda_L$  described above, the adaptive LASSO has the following expression

$$\widehat{\beta}^{\text{adapt}} := \underset{\beta \in \mathbb{R}^p}{\text{argmin}} \frac{1}{2} \|Y - X\beta\|_2^2 + \lambda_L \sum_{i=1}^p w(\beta_i) |\beta_i|. \quad (7)$$

In all our simulations LASSO is calculated with *glmnet*.

### 5.3 Numerical comparisons

The rows of the design matrix  $X$  are sampled as the independent vectors from the multivariate Gaussian distribution, as in setting 1 and 2. All numerical experiments are performed with a particular observation of  $X$  (the same as the one used in the previous subsection). We set  $\beta^0 \in \mathbb{R}^p$  such that  $k := \text{card}\{i \mid \beta_i^0 \neq 0\}$  with  $k = \{5, 20\}$ ,  $\{i \mid \beta_i^0 \neq 0\}$  is a  $k$  sample without replacement of  $\{1, \dots, p\}$ . The non null components of  $\beta^0$  have a uniform distribution  $\{-t, t\}$  where  $t > 0$ . Additionally, for strongly and positively correlated explanatory variables we consider the set-up where all non-zero coefficients are equal to  $t$ . In all simulations the error term is generated as  $\varepsilon \sim \mathcal{N}(0, Id_n)$ .

Figures 4-6 provide the comparison between the following sign estimators.

- The sign estimator **L** is derived from LASSO with  $\lambda = \lambda_L$ .
- The sign estimator **adL** is derived from the adaptive LASSO estimator, described in (7).
- The sign estimator **BPS** is derived from the thresholded BP, with threshold selected as in [12].
- The sign estimator **BPkn** is derived from the thresholded BP, with a threshold given by the "knockoff" methodology described above.
- The sign estimator **Lkn** is derived from the thresholded LASSO with  $\lambda = \lambda_{AMP}$  and with a threshold given by the "knockoff" methodology described above.
- The sign estimator **Lkns** is derived from the thresholded LASSO with  $\lambda = 0.5\lambda_{AMP}$  and with a threshold given by the "knockoff" methodology described above.

In order to recover the sign of  $\beta^0$ , null components of  $\beta^0$  have to be estimated simultaneously at zero. This naive remark motivate us to report the curves illustrating the following statistical properties as the function of the signal strength:

- **FWER** is the proportion of 1000 replicates that at least one null components of  $\beta^0$  be not estimated at zero.

We report the curve illustrating the probability to recover the sign as the function of the signal strenght:

- **Probability** is the proportion of 1000 replicates for which the sign is recovered.

Figure 6-8 illustrate that the upper bound for the probability of LASSO sign estimator is reached and the FWER is controlled when non null component of  $\beta^0$  are large (*i.e* when  $t$  is large). On the other hand, thresholded LASSO and thresholded BP can appropriately identify  $S(\beta^0)$  when the identifiability condition holds. Indeed, when  $k \in \{5, 20\}$  as illustrated in figures 4 and 5, the identifiability condition occurs and thus sign estimators derived from thresholded LASSO and thresholded BP recover  $S(\beta^0)$  as soon as the threshold is

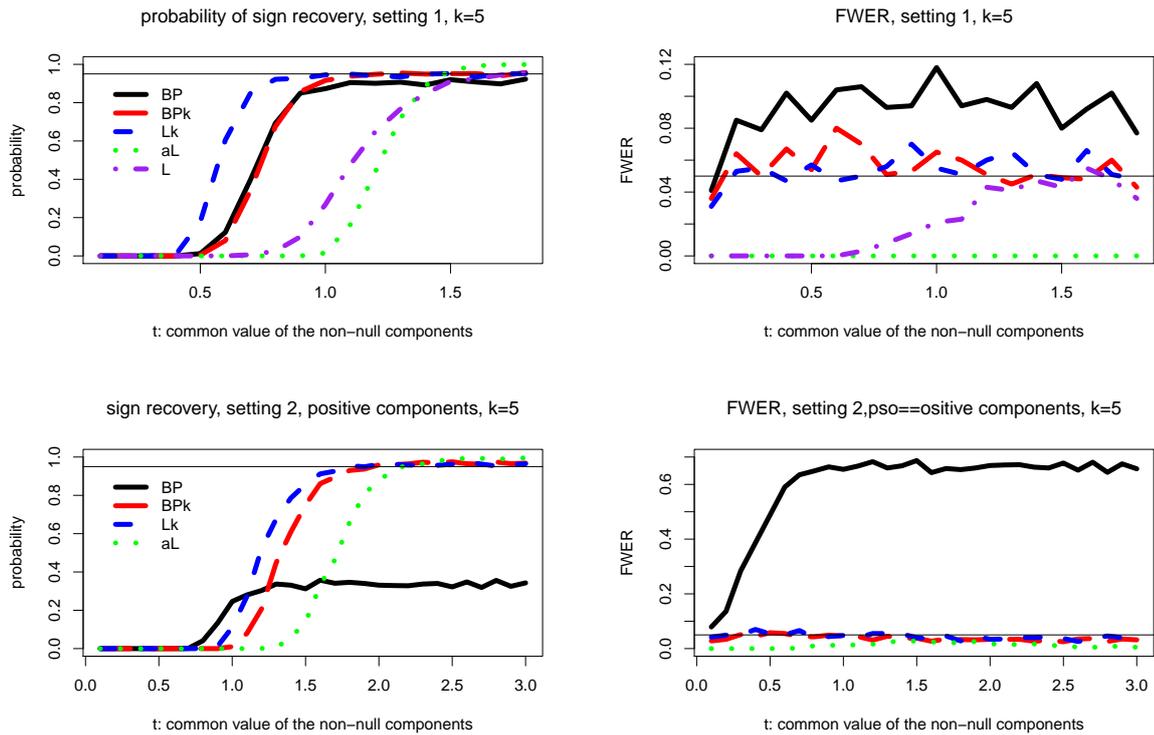


Figure 6: This figure provides the FWER and the probability to recover  $S(\beta^0)$  for each sign estimator and when  $X$  is the design matrix given in setting 1. Figures on the left provide the probability to recover  $S(\beta^0)$  (on the y-axis) as a function of  $t$ , where  $t$  measures how large are the non-null components of  $\beta^0$ . Figures on the right provide the FWER (on the y-axis) as a function of  $t$  (on the x-axis). Among these sign estimators, one notices that the thresholded LASSO sign estimator is the one which recovers  $S(\beta^0)$  with the largest probability. These sign estimators recover approximately  $S(\beta^0)$  with a probability close to 0.95 when  $t$  is large.

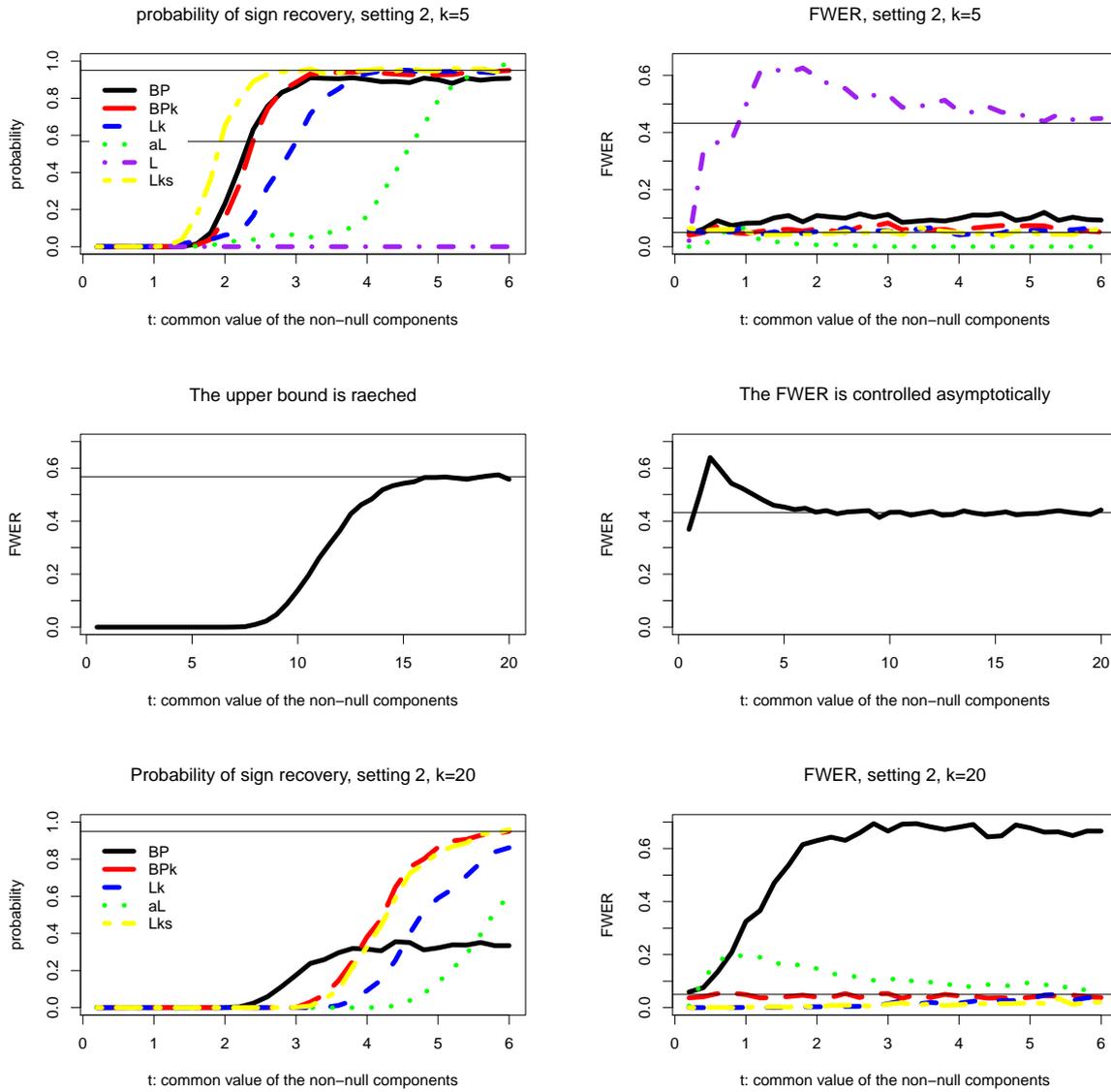


Figure 7: This figure provides the FWER and the probability to recover  $S(\beta^0)$  for each sign estimators and when  $X$  is the design matrix given in setting 2. Figures on the left provide the probability to recover  $S(\beta^0)$  (on the y-axis) as a function of  $t$ , where  $t$  measures how large are the non-null components of  $\beta^0$ . Figures on the right provide the FWER (on the y-axis) as a function of  $t$  (on the x-axis). The horizontal lines  $y = 0.55$  and  $y = 0.45$  are the average value upper bound for the probability LASSO sign recovery and average value for the FWER associated to LASSO given by the proposition 1. One notices that the upper-bound is approximately reached and the FWER is approximately controlled when  $t$  is very large as illustrated by figures in the middle. Sign estimators (except LASSO sign estimator) recover approximately  $S(\beta^0)$  with a probability close to 0.95 when  $t$  is large.

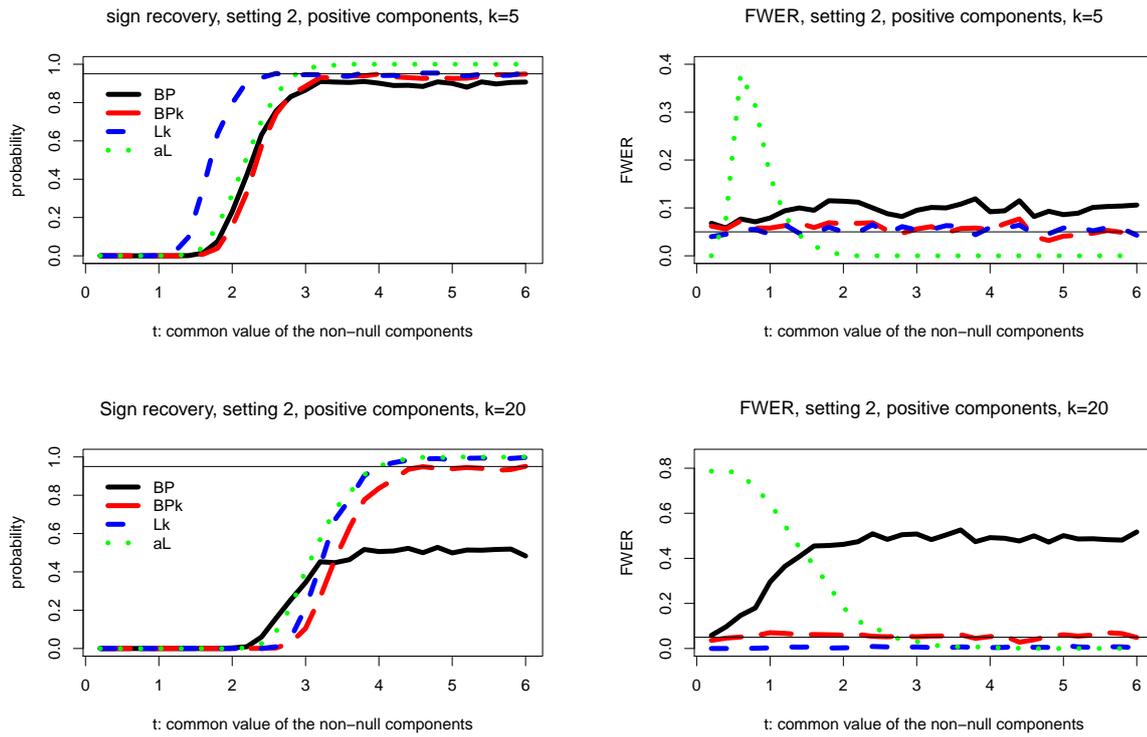


Figure 8: This figure provides the FWER and the probability to recover  $S(\beta^0)$  for each sign estimator and when  $X$  is the design matrix given in setting 2 and non-null components of  $\beta^0$  are positive. Figures on the left provide the probability to recover  $S(\beta^0)$  (on the y-axis) as a function of  $t$ , where  $t$  measures how large are the non-null components of  $\beta^0$ . Figures on the right provide the FWER (on the y-axis) as a function of  $t$  (on the x-axis). These sign estimators recover approximately  $S(\beta^0)$  with a probability close to 0.95 when  $t$  is large.

well calibrated and the non null components are large enough. In our simulated setup, thresholded BP performs pretty well but is never optimal. Indeed using an appropriate tuning parameter  $\lambda$ , the probability to recover  $S(\beta^0)$  is larger with thresholded LASSO than with thresholded BP. When entries of  $X$  are i.i.d  $\mathcal{N}(0, 1)$ , the optimal value of  $\lambda$  selected by AMP theory provides a thresholded LASSO for which the derived sign estimator is the best one to recover  $S(\beta^0)$ . One notices that the threshold selection provided in Descloux and Sardy [12] does not allow to recover  $S(\beta^0)$  with a large probability when  $\beta^0$  has lot of large components (intuitively when  $\beta^0$  is far from  $\mathbf{0}$ ). Instead, our heuristic application of the knockoff methodology allows for almost perfect control of FWER at level 0.05. Consequently, when non null components of  $\beta^0$  are large enough and when the threshold is given by knockoff methodology, sign estimator derived from thresholded LASSO (resp. thresholded BP) recovers  $S(\beta^0)$  with a probability close to 0.95.

## 6 Conclusion

This article main focus on theoretical properties of sign estimators derived from LASSO, thresholded LASSO and thresholded BPDN. We provided an upper bound for LASSO sign recovery; this upper bound is reached when  $N$  tends to  $+\infty$  and the identifiability condition holds. In addition, when the irrepresentable condition occurs (implying that the identifiability condition occurs) and when  $N$  is large, we have shown that  $\lambda$  can be selected appropriately in order to control the FWER at an arbitrary level.

When  $S(\beta^*)$  is identifiable with respect to the  $l^1$  norm and  $N$  is large enough, we have shown that sign estimators derived from thresholded LASSO and thresholded BPDN recover  $S(\beta^*)$ . On the other hand, if  $S(\beta^*)$  is not identifiable with respect to the  $l^1$  norm, sign estimators derived from thresholded LASSO and thresholded cannot recover  $S(\beta^*)$  even if  $N$  is extremely large.

We have introduced identifiability curve (resp. irrepresentability curve) which is useful to know for which sparsity  $\beta^*$  is identifiable with respect to the  $l^1$  norm (resp. for which sparsity  $\beta^*$  the irrepresentable condition holds).

The performances of sign estimators derived from LASSO, thresholded LASSO and thresholded BPDN depend obviously from the tuning parameter, the regularization parameter and the threshold. We have illustrated that AMP theory and knockoff methodology are useful to select these two parameters. Our simulations show that thresholded LASSO and thresholded BPDN sign estimators outperform adaptive LASSO and LASSO sign estimators.

## Acknowledgments

We would like to thank Emmanuel J. Candès for helpful comments. The research of Małgorzata Bogdan was funded by the NCN grant 2016/23/B/ST1/00454.

## 7 appendix

### 7.1 Proof of the theorem 1

First, let us provide a lemma which is useful to prove both theorems 1 and 2. This lemma partially proved the theorem 1. Indeed, according to the lemma 1, when  $\beta^*$  satisfies the identifiability condition then the following asymptotic result holds

$$\lim_{N \rightarrow +\infty} \mathbb{P}(\forall i \in I, S(\hat{\beta}_i^{\text{lasso}}(\lambda, N)) = S(\beta_i^*)) = 1.$$

**Lemma 1** *Let us remind that  $\tilde{\beta} \in \mathbb{R}^p$  is the unique solution of the problem: minimize  $\|\beta\|_1$  subject to  $X\beta =$*

$X\beta^*$  and let  $\varepsilon \in \mathbb{R}^n$  be a fixed vector. Then, the following asymptotic result holds

$$\lim_{N \rightarrow +\infty} \widehat{\beta}^{\text{lasso}}(\varepsilon, N)/N = \tilde{\beta}.$$

**Proof:** Because  $\widehat{\beta}^{\text{lasso}}(\varepsilon, N)$  is the LASSO estimator as defined in (2) then the following inequality occurs

$$\frac{1}{2} \|Y - X\widehat{\beta}^{\text{lasso}}(\varepsilon, N)\|_2^2 + \lambda \|\widehat{\beta}^{\text{lasso}}(\varepsilon, N)\|_1 \leq \frac{1}{2} \|Y - X(N\tilde{\beta})\|_2^2 + \lambda \|N\tilde{\beta}\|_1.$$

Since  $Y - X(N\tilde{\beta}) = \varepsilon$  one may deduce the following inequalities

$$\begin{aligned} \lambda \|\widehat{\beta}^{\text{lasso}}(\varepsilon, N)\|_1 &\leq \frac{1}{2} \|\varepsilon\|_2^2 + \lambda \|N\tilde{\beta}\|_1, \\ \Rightarrow \|\widehat{\beta}^{\text{lasso}}(\varepsilon, N)/N\|_1 &\leq \frac{\|\varepsilon\|_2^2}{2\lambda N} + \|\tilde{\beta}\|_1. \end{aligned} \quad (8)$$

In addition, Cauchy-Schwarz inequality gives the following implications

$$\begin{aligned} &\frac{1}{2} \|\varepsilon + X(N\tilde{\beta}) - X\widehat{\beta}^{\text{lasso}}(\varepsilon, N)\|_2^2 + \lambda \|\widehat{\beta}^{\text{lasso}}(\varepsilon, N)\|_1 \leq \frac{1}{2} \|\varepsilon\|_2^2 + \lambda \|N\tilde{\beta}\|_1, \\ \Rightarrow -\|\varepsilon\|_2 \|X(N\tilde{\beta}) - X\widehat{\beta}^{\text{lasso}}(\varepsilon, N)\|_2 + \frac{1}{2} \|X(N\tilde{\beta}) - X\widehat{\beta}^{\text{lasso}}(\varepsilon, N)\|_2^2 + \lambda \|\widehat{\beta}^{\text{lasso}}(\varepsilon, N)\|_1 &\leq \lambda \|N\tilde{\beta}\|_1, \\ \Rightarrow -\frac{\|\varepsilon\|_2}{N} \left\| X \left( \tilde{\beta} - \frac{\widehat{\beta}^{\text{lasso}}(\varepsilon, N)}{N} \right) \right\|_2 + \frac{1}{2} \left\| X \left( \tilde{\beta} - \frac{\widehat{\beta}^{\text{lasso}}(\varepsilon, N)}{N} \right) \right\|_2^2 + \frac{\lambda}{N} \left\| \frac{\widehat{\beta}^{\text{lasso}}(\varepsilon, N)}{N} \right\|_1 &\leq \frac{\lambda}{N} \|\tilde{\beta}\|_1 \end{aligned} \quad (9)$$

To prove that  $\lim_{N \rightarrow +\infty} \widehat{\beta}^{\text{lasso}}(\varepsilon, N)/N = \tilde{\beta}$  is enough to prove that for an arbitrary increasing sequence  $(N_r)_{r \in \mathbb{N}^*}$  such that  $\lim_{r \rightarrow +\infty} N_r = +\infty$  then  $\lim_{r \rightarrow +\infty} \widehat{\beta}^{\text{lasso}}(\varepsilon, N_r)/N_r = \tilde{\beta}$ . Let us notice that, according to (8), the sequence  $(\widehat{\beta}^{\text{lasso}}(\varepsilon, N_r)/N_r)_{r \in \mathbb{N}^*}$  is bounded (by  $1 + \|\varepsilon\|_2^2/(\lambda N_0)$ ). Consequently, to prove that  $\lim_{r \rightarrow +\infty} (\widehat{\beta}^{\text{lasso}}(\varepsilon, N_r)/N_r) = \tilde{\beta}$  it is sufficient to show that  $\tilde{\beta}$  is the unique limit point of this sequence. Let  $(\widehat{\beta}^{\text{lasso}}(\varepsilon, N_{\phi(r)})/N_{\phi(r)})_{r \in \mathbb{N}^*}$  be a converging subsequence to  $l$  (with  $\phi : \mathbb{N}^* \rightarrow \mathbb{N}^*$  strictly increasing). By (8) and (9) one deduces that

$$X\tilde{\beta} = Xl \text{ and } \|l\|_1 \leq \|\tilde{\beta}\|_1.$$

By construction,  $\tilde{\beta}$  is identifiable up to its  $l^1$  norm (as the unique solution of a basis pursuit problem). Therefore, one may deduce that  $\tilde{\beta} = l$  thus  $\tilde{\beta}$  is the unique limit point. Consequently,  $\lim_{r \rightarrow +\infty} \widehat{\beta}^{\text{lasso}}(\varepsilon, N_r)/N_r = \tilde{\beta}$ , which implies that

$$\lim_{N \rightarrow +\infty} \widehat{\beta}^{\text{lasso}}(\varepsilon, N)/N = \tilde{\beta}.$$

□

**Proof of the theorem 1:** Let  $A$  be the set  $A := \text{supp}(\widehat{\beta}^{\text{lasso}}(\lambda))$ .

**Upper bound**) Let us give two expressions met by the LASSO estimator as defined in (2). The vector  $\widehat{\beta}^{\text{lasso}}(\lambda)$  is the LASSO estimator if and only if the following two inequalities occur simultaneously.

$$X'_A(Y - X\widehat{\beta}^{\text{lasso}}(\lambda)) = \lambda S(\widehat{\beta}_A^{\text{lasso}}(\lambda)), \quad (10)$$

$$\|X'_A(Y - X\widehat{\beta}^{\text{lasso}}(\lambda))\|_\infty \leq \lambda. \quad (11)$$

These two expressions are given in Bühlmann and van de Geer [5] page 15 or in the proof of the theorem 1 of Zou [32]. Using the equality (10) and the inequality (11), we are going to show that if  $S(\widehat{\beta}^{\text{lasso}}(\lambda)) = S(\beta^*)$  then the following event holds

$$\left\| X'_I X_I (X'_I X_I)^{-1} S(\beta_I^*) + \frac{1}{\lambda} X'_I (Id - X_I (X'_I X_I)^{-1} X'_I) \varepsilon \right\|_\infty \leq 1.$$

Let us assume that  $S(\widehat{\beta}^{\text{lasso}}(\lambda)) = S(\beta^*)$  thus  $A = I$  (where  $I = \text{supp}(\beta^*)$ ). Since  $Y = X\beta^* + \varepsilon = X_I\beta_I^* + \varepsilon$  and  $X\widehat{\beta}^{\text{lasso}}(\lambda) = X_I\widehat{\beta}_I^{\text{lasso}}(\lambda)$  then the equality (10) and the inequality (11) lead to the following expressions

$$X'_I \left( \varepsilon + X_I(\beta_I^* - \widehat{\beta}_I^{\text{lasso}}(\lambda)) \right) = \lambda S(\beta_I^*), \quad (12)$$

$$\left\| X'_I \left( \varepsilon + X_I(\beta_I^* - \widehat{\beta}_I^{\text{lasso}}(\lambda)) \right) \right\|_\infty \leq \lambda. \quad (13)$$

The equality (12) assures that

$$\beta_I^* - \widehat{\beta}_I^{\text{lasso}}(\lambda) = (X'_I X_I)^{-1} (\lambda S(\beta_I^*) - X'_I \varepsilon).$$

Let us notice that since  $\ker(X_I) = \mathbf{0}$  then the Gram matrix  $X'_I X_I$  is invertible. Using the previous expression in the inequality (13) gives

$$\begin{aligned} \|X'_I X_I (X'_I X_I)^{-1} (\lambda S(\beta_I^*) - X'_I \varepsilon) + X'_I \varepsilon\|_\infty &\leq \lambda, \\ \left\| X'_I X_I (X'_I X_I)^{-1} S(\beta_I^*) + \frac{1}{\lambda} X'_I (Id - X_I (X'_I X_I)^{-1} X'_I) \varepsilon \right\|_\infty &\leq 1. \end{aligned}$$

Consequently, one deduces the following inequality

$$\mathbb{P} \left( S(\widehat{\beta}^{\text{lasso}}(\lambda)) = S(\beta^*) \right) \leq \underbrace{\mathbb{P} \left( \left\| X'_I X_I (X'_I X_I)^{-1} S(\beta_I^*) + \frac{1}{\lambda} X'_I (Id - X_I (X'_I X_I)^{-1} X'_I) \varepsilon \right\|_\infty \leq 1 \right)}_{=\mathbb{P}(\|\zeta_{X, \lambda, S(\beta^*)}\|_\infty \leq 1) = \bar{\gamma}}.$$

**Sharpness of the upper bound**) Since the upper bound does not depend from  $N > 0$  then

$$\limsup_{N \rightarrow +\infty} \mathbb{P} \left( S(\widehat{\beta}^{\text{lasso}}(\lambda, N)) = S(\beta^*) \right) \leq \bar{\gamma}.$$

So it remains to prove that  $\liminf_{N \rightarrow +\infty} \mathbb{P} \left( S(\widehat{\beta}^{\text{lasso}}(\lambda, N)) = S(\beta^*) \right) \geq \mathbb{P} \left( \|\zeta_{X, \lambda, S(\beta^*)}\|_{\infty} < 1 \right) = \gamma$ .

Let us assume that the following events hold simultaneously

$$X_I'(Y - X\widehat{\beta}^{\text{lasso}}(\lambda)) = \lambda S(\beta_I^*) \text{ and } \underbrace{\left\| X_I'X_I(X_I'X_I)^{-1}\lambda S(\beta_I^*) + X_I'(Id - X_I(X_I'X_I)^{-1}X_I')\varepsilon \right\|_{\infty}}_{=\|\zeta_{X, \lambda, S(\beta^*)}\|_{\infty} < 1} < \lambda. \quad (14)$$

We aim to show that the inequalities given above imply that  $\widehat{\beta}^{\text{lasso}}(\lambda) = \mathbf{0}$ . For convenience, let us set  $H$  be the projection matrix  $H := X_I(X_I'X_I)^{-1}X_I'$ . When (14) occurs then the following inequalities holds

$$\begin{aligned} \left\| X_I'H(Y - X\widehat{\beta}^{\text{lasso}}(\lambda)) + X_I'(Id - H)\varepsilon \right\|_{\infty} &< \lambda, \\ \left\| X_I'(H(Y - X\widehat{\beta}^{\text{lasso}}(\lambda)) + (Id - H)\varepsilon) \right\|_{\infty} &< \lambda, \\ \left\| X_I'(Y - X\widehat{\beta}^{\text{lasso}}(\lambda) + X_I\widehat{\beta}_I^{\text{lasso}}(\lambda) + HX_I\widehat{\beta}_I^{\text{lasso}}(\lambda)) \right\|_{\infty} &< \lambda. \end{aligned} \quad (15)$$

The inequality (15) comes from the following two identities

$$\begin{aligned} HY &= H(X(N\beta^*)) + H\varepsilon = H(X_I(N\beta_I^*)) + H\varepsilon = X_I(N\beta_I^*) + H\varepsilon = X(N\beta^*) + H\varepsilon \text{ and,} \\ HX\widehat{\beta}^{\text{lasso}}(\lambda) &= HX_I\widehat{\beta}_I^{\text{lasso}}(\lambda) + HX_I\widehat{\beta}_I^{\text{lasso}}(\lambda) = X\widehat{\beta}^{\text{lasso}}(\lambda) - X_I\widehat{\beta}_I^{\text{lasso}}(\lambda) + HX_I\widehat{\beta}_I^{\text{lasso}}(\lambda). \end{aligned}$$

Let  $v$  be the vector  $v := X_I'(Y - X\widehat{\beta}^{\text{lasso}}(\lambda) + X_I\widehat{\beta}_I^{\text{lasso}}(\lambda) + HX_I\widehat{\beta}_I^{\text{lasso}}(\lambda))$ , we are going to see that the inequality (15) implies that  $\widehat{\beta}_I^{\text{lasso}}(\lambda) = \mathbf{0}$ . Let us assume that  $\widehat{\beta}_I^{\text{lasso}}(\lambda) \neq \mathbf{0}$  then, on the first hand, the following inequality occurs

$$\widehat{\beta}_I^{\text{lasso}}(\lambda)'v \leq \|\widehat{\beta}_I^{\text{lasso}}(\lambda)\|_1 \|v\|_{\infty} < \lambda \|\widehat{\beta}_I^{\text{lasso}}(\lambda)\|_1. \quad (16)$$

According to (10) the identity  $\widehat{\beta}_i^{\text{lasso}}(\lambda)X_i'(Y - X\widehat{\beta}_i^{\text{lasso}}(\lambda)) = \lambda|\widehat{\beta}_i^{\text{lasso}}(\lambda)|$  occurs. Consequently, on the other hand, the following inequalities hold

$$\begin{aligned} \widehat{\beta}_I^{\text{lasso}}(\lambda)'v &= \widehat{\beta}_I^{\text{lasso}}(\lambda)'X_I'(Y - X\widehat{\beta}^{\text{lasso}}(\lambda) + X_I\widehat{\beta}_I^{\text{lasso}}(\lambda) - HX_I\widehat{\beta}_I^{\text{lasso}}(\lambda)), \\ &= \lambda\|\widehat{\beta}_I^{\text{lasso}}(\lambda)\|_1 + \widehat{\beta}_I^{\text{lasso}}(\lambda)'X_I'(Id - H)X_I\widehat{\beta}_I^{\text{lasso}}(\lambda), \\ &\geq \lambda\|\widehat{\beta}_I^{\text{lasso}}(\lambda)\|_1. \end{aligned} \quad (17)$$

The last inequality occurs because the projection matrix  $Id - H$  is positive semi-definite. Inequalities (16) and (17) provide a contradiction which implies that  $\widehat{\beta}_I^{\text{lasso}}(\lambda) = \mathbf{0}$ .

According to (10), the following implication holds

$$S(\widehat{\beta}_I^{\text{lasso}}(\lambda, N)) = S(\beta_I^*) \Rightarrow X_I'(Y - X\widehat{\beta}^{\text{lasso}}(\lambda, N)) = \lambda S(\beta_I^*).$$

Because  $\beta^*$  is identifiable with respect to the  $l^1$  norm then, according to the lemma 1, the following convergence in probability occurs

$$\lim_{N \rightarrow +\infty} \mathbb{P}(S(\widehat{\beta}_I^{\text{lasso}}(\lambda, N)) = S(\beta_I^*)) = \lim_{N \rightarrow +\infty} \mathbb{P}(X_I'(Y - X\widehat{\beta}^{\text{lasso}}(\lambda, N)) = \lambda S(\beta_I^*)) = 1. \quad (18)$$

Using this asymptotic result and since when (14) occurs then  $\widehat{\beta}_I^{\text{lasso}}(\lambda, N) = \mathbf{0}$ , one may deduce the following inequalities

$$\begin{aligned} \liminf_{N \rightarrow +\infty} \mathbb{P}\left(S(\widehat{\beta}^{\text{lasso}}(\lambda, N)) = S(\beta^*)\right) &= \liminf_{N \rightarrow +\infty} \mathbb{P}\left(S(\widehat{\beta}_I^{\text{lasso}}(\lambda, N)) = S(\beta_I^*) \text{ and } \widehat{\beta}_I^{\text{lasso}}(\lambda, N) = \mathbf{0}\right), \\ &= \liminf_{N \rightarrow +\infty} \mathbb{P}(\widehat{\beta}_I^{\text{lasso}}(\lambda, N) = \mathbf{0}), \\ &\geq \liminf_{N \rightarrow +\infty} \mathbb{P}\left(X_I'(Y - X\widehat{\beta}^{\text{lasso}}(\lambda, N)) = S(\beta_I^*) \text{ and } \|\zeta_{X, \lambda, S(\beta^*)}\|_\infty < 1\right), \\ &\geq \liminf_{N \rightarrow +\infty} \mathbb{P}(\|\zeta_{X, \lambda, S(\beta^*)}\|_\infty < 1) = \gamma. \end{aligned}$$

**Asymptotic full power and asymptotic control of the FWER)** According to (18), asymptotically the power is equal to 1, namely  $\lim_{N \rightarrow +\infty} \mathbb{P}(\forall i \in I, S(\widehat{\beta}_i^{\text{lasso}}(\lambda, N)) = S(\beta_i^*)) = 1$ . Now let us prove that the FWER is controlled asymptotically. Using asymptotic results given above one may deduce the following inequalities.

$$\begin{aligned} \bar{\gamma} &\geq \limsup_{N \rightarrow +\infty} \mathbb{P}(S(\widehat{\beta}(\lambda, N)) = S(\beta^*)), \\ &\geq \limsup_{N \rightarrow +\infty} \mathbb{P}(\forall i \in I, S(\widehat{\beta}_i(\lambda, N)) = S(\beta_i^*) \text{ and } \forall i \notin I, \widehat{\beta}_i(\lambda, N) = 0), \\ &\geq \limsup_{N \rightarrow +\infty} \mathbb{P}(\forall i \notin I, \widehat{\beta}_i(\lambda, N) = 0). \end{aligned} \quad (19)$$

The last inequality come from (18). Similarly, we have

$$\gamma \leq \liminf_{N \rightarrow +\infty} \mathbb{P}(\forall i \notin I, \widehat{\beta}_i(\lambda, N) = 0). \quad (20)$$

Consequently, by taking the complement to 1 of the inequalities given in (19) and (20), one may deduce that

$$\liminf_{N \rightarrow +\infty} \mathbb{P}(\exists i \notin I, \widehat{\beta}_i(\lambda, N) \neq 0) \geq 1 - \bar{\gamma} \text{ and } \limsup_{N \rightarrow +\infty} \mathbb{P}(\exists i \notin I, \widehat{\beta}_i(\lambda, N) \neq 0) \leq 1 - \gamma.$$

□

## Proof of the theorem 2

The lemma 2 provides the same result for BPDN than the lemma 1 for LASSO. These both lemmas are the keystone to prove the theorem 2. The proof of lemma 2 is partially inspired from the one given in Candès et al. [9].

**Lemma 2** Let us remind that  $\tilde{\beta} \in \mathbb{R}^p$  is the unique solution of the problem: minimize  $\|\beta\|_1$  subject to  $X\beta = X\beta^*$  and let  $\varepsilon \in \mathbb{R}^n$  be a fixed vector. Then, the following asymptotic result holds

$$\lim_{N \rightarrow +\infty} \widehat{\beta}^{\text{bpdn}}(\varepsilon, N)/N = \tilde{\beta}.$$

**Proof:** Let us define  $u(\varepsilon) \in \mathbb{R}^p$  as follows

$$u(\varepsilon) := \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \|\beta\|_1 \text{ subject to } X\beta = \varepsilon.$$

Because  $X(N\beta^*) = X(N\tilde{\beta})$  and  $X(u(\varepsilon)) = \varepsilon$ , we have  $Y(\varepsilon) = X(N\tilde{\beta} + u(\varepsilon))$  and because  $\widehat{\beta}^{\text{bpdn}}(\varepsilon, N)$  is an admissible point of (3) one deduces the following inequality

$$\left\| \frac{1}{N} X \widehat{\beta}^{\text{bpdn}}(\varepsilon, N) - X \tilde{\beta} \right\|_2 \leq \left\| \frac{1}{N} X \widehat{\beta}^{\text{bpdn}}(\varepsilon, N) - \frac{1}{N} Y \right\|_2 + \left\| \frac{1}{N} Y - X \tilde{\beta} \right\|_2 \leq \frac{\sqrt{R}}{N} + \frac{\|Xu(\varepsilon)\|_2}{N}. \quad (21)$$

Because  $N\tilde{\beta} + u(\varepsilon)$  is an admissible point of the problem (3) and because  $\widehat{\beta}^{\text{bpdn}}(\varepsilon, N)$  is the minimizer of (3), one deduces the following inequalities hold

$$\frac{1}{N} \|\widehat{\beta}^{\text{bpdn}}(\varepsilon, N)\|_1 \leq \frac{1}{N} \|N\tilde{\beta} + u(\varepsilon)\|_1 \leq \|\tilde{\beta}\|_1 + \frac{\|u(\varepsilon)\|_1}{N}. \quad (22)$$

To prove that  $\lim_{N \rightarrow +\infty} \widehat{\beta}^{\text{bpdn}}(\varepsilon, N)/N = \tilde{\beta}$  is enough to prove that for an arbitrary increasing sequence  $(N_r)_{r \in \mathbb{N}^*}$  such that  $\lim_{r \rightarrow +\infty} N_r = +\infty$  then  $\lim_{r \rightarrow +\infty} \widehat{\beta}^{\text{bpdn}}(\varepsilon, N_r)/N_r = \tilde{\beta}$ . Let us notice that the sequence  $(\widehat{\beta}^{\text{bpdn}}(\varepsilon, N_r)/N_r)_{r \in \mathbb{N}^*}$  is bounded (by  $1 + \|u(\varepsilon)\|_1/N_0$ ). Consequently, to prove that  $\lim_{r \rightarrow +\infty} (\widehat{\beta}^{\text{bpdn}}(\varepsilon, N_r)/N_r) = \tilde{\beta}$  it is sufficient to show that  $\tilde{\beta}$  is the unique limit point of this sequence. Let  $(\widehat{\beta}^{\text{bpdn}}(\varepsilon, N_{\phi(r)})/N_{\phi(r)})_{r \in \mathbb{N}^*}$  be a converging subsequence to  $l$  (with  $\phi : \mathbb{N}^* \rightarrow \mathbb{N}^*$  strictly increasing). By (21) and (22) one deduces that

$$X\tilde{\beta} = Xl \text{ and } \|l\|_1 \leq \|\tilde{\beta}\|_1.$$

By construction of  $\tilde{\beta}$  (as a unique solution of a basis pursuit problem), one deduces that  $\tilde{\beta} = l$  thus  $\tilde{\beta}$  is the unique limit point. Consequently,  $\lim_{r \rightarrow +\infty} \widehat{\beta}^{\text{bpdn}}(\varepsilon, N_r)/N_r = \tilde{\beta}$ , which implies that

$$\lim_{N \rightarrow +\infty} \widehat{\beta}^{\text{bpdn}}(\varepsilon, N)/N = \tilde{\beta}.$$

□

**Proof of the theorem 2:**

**Separation property:** Let us set  $\eta_0 > 0$  such that  $\eta_0 < \min\{|\tilde{\beta}_i|, i \in \mathcal{B}\}/2$ . The lemmas 1 and 2 prove the

convergence of  $\widehat{\beta}(\varepsilon, N)/N$  to  $\tilde{\beta}$  when  $N$  tends to  $+\infty$ . Consequently, there exists  $N_0(\varepsilon) \in (0, +\infty)$  such that

$$\forall N \geq N_0(\varepsilon), \|\widehat{\beta}(\varepsilon, N)/N - \tilde{\beta}\|_\infty \leq \eta_0 \Leftrightarrow \forall N \geq N_0(\varepsilon), \forall i \in \{1, \dots, p\}, \left| \widehat{\beta}_i(\varepsilon, N)/N - \tilde{\beta}_i \right| \leq \eta_0.$$

Consequently, when  $N \geq N_0(\varepsilon)$ , whatever  $i \notin \mathcal{B}$  (thus when  $\tilde{\beta}_i = 0$ ) the following inequalities hold

$$\begin{aligned} & \forall i \notin \mathcal{B}, \left| \widehat{\beta}_i(\varepsilon, N)/N \right| \leq \eta_0, \\ \Rightarrow & -N\eta_0 \leq \min_{i \notin \mathcal{B}} \left\{ \widehat{\beta}_i(\varepsilon, N) \right\} \leq \max_{i \notin \mathcal{B}} \left\{ \widehat{\beta}_i(\varepsilon, N) \right\} \leq N\eta_0. \end{aligned}$$

Whatever  $i \in \mathcal{B}^+$  (thus when  $\tilde{\beta}_i > 0$ ) the following inequalities hold

$$\begin{aligned} & \forall i \in \mathcal{B}^+, \widehat{\beta}_i(\varepsilon, N)/N \geq -\left| \widehat{\beta}_i(\varepsilon, N)/N - \tilde{\beta}_i \right| + \tilde{\beta}_i, \\ \Rightarrow & \min_{i \in \mathcal{B}^+} \left\{ \widehat{\beta}_i(\varepsilon, N)/N \right\} \geq -\eta_0 + \min\{|\tilde{\beta}_i|, i \in \mathcal{B}\} > \eta_0, \\ \Rightarrow & \min_{i \in \mathcal{B}^+} \left\{ \widehat{\beta}_i(\varepsilon, N) \right\} > N\eta_0. \end{aligned}$$

Whatever  $i \in \mathcal{B}^-$  (thus when  $\tilde{\beta}_i < 0$ ) the following inequalities hold

$$\begin{aligned} & \forall i \in \mathcal{B}^-, \widehat{\beta}_i(\varepsilon, N)/N \leq \left| \widehat{\beta}_i(\varepsilon, N)/N - \tilde{\beta}_i \right| + \tilde{\beta}_i, \\ \Rightarrow & \max_{i \in \mathcal{B}^-} \left\{ \widehat{\beta}_i(\varepsilon, N)/N \right\} \leq \eta_0 - \min\{|\tilde{\beta}_i|, i \in \mathcal{B}\} < -\eta_0, \\ \Rightarrow & \max_{i \in \mathcal{B}^-} \left\{ \widehat{\beta}_i(\varepsilon, N) \right\} < -N\eta_0. \end{aligned}$$

Finally, when  $N \geq N_0(\varepsilon)$  then

$$\max_{i \in \mathcal{B}^-} \left\{ \widehat{\beta}_i(\varepsilon, N) \right\} < \min_{i \notin \mathcal{B}} \left\{ \widehat{\beta}_i(\varepsilon, N) \right\} \leq \max_{i \notin \mathcal{B}} \left\{ \widehat{\beta}_i(\varepsilon, N) \right\} < \min_{i \in \mathcal{B}^+} \left\{ \widehat{\beta}_i(\varepsilon, N) \right\}.$$

**Sign recovery:** If  $\beta^*$  is identifiable with respect to the  $l^1$  norm then  $\beta^* = \tilde{\beta}$  and consequently,  $S(\tilde{\beta}) = S(\beta^*)$ .

Reciprocally, let us assume that  $S(\tilde{\beta}) = S(\beta^*)$ . Because, by construction,  $\tilde{\beta}$  is identifiable with respect to the  $l^1$  norm and because  $S(\tilde{\beta}) = S(\beta^*)$  then, according to the proposition 2,  $\beta^*$  is identifiable with respect to the  $l^1$  norm.  $\square$

## Proof of propositions

**Proof of the proposition 2:** According to Daubechies et al. [11],  $\beta^*$  is identifiable with respect to the  $l^1$

norm if and only if the following inequality holds

$$\forall h \in \ker(X) \setminus \{\mathbf{0}\}, \left| \sum_{i \in \text{supp}(\beta^*)} S(\beta_i^*) h_i \right| < \sum_{i \notin \text{supp}(\beta^*)} |h_i|.$$

Because  $S(\tilde{\beta}) = S(\beta^*)$  implies  $\text{supp}(\tilde{\beta}) = \text{supp}(\beta^*)$ , the following inequality holds

$$\forall h \in \ker(X) \setminus \{\mathbf{0}\}, \left| \sum_{i \in \text{supp}(\tilde{\beta})} S(\tilde{\beta}_i) h_i \right| < \sum_{i \notin \text{supp}(\tilde{\beta})} |h_i|.$$

Consequently, the parameter  $\tilde{\beta}$  is identifiable with respect to the  $l^1$  norm.  $\square$

**Proof of the proposition 1:** From Daubechies et al. [11],  $\beta^*$  is a parameter having a minimal  $l^1$  norm, namely  $X\beta^* = X\gamma \Rightarrow \|\gamma\|_1 \geq \|\beta^*\|_1$  holds, if and only if the following inequality occurs

$$\forall h \in \ker(X), \left| \sum_{i \in \mathcal{A}} S(\beta_i^*) h_i \right| \leq \sum_{i \notin \mathcal{A}} |h_i|. \quad (23)$$

We are going to show that when the irrerepresentable condition holds for  $\beta^*$  then the inequality (23) holds.

Let  $h \in \ker(X)$  and let us remind that  $h_{\mathcal{A}}$  and  $h_{\bar{\mathcal{A}}}$  denote respectively vectors  $(h_i)_{i \in \mathcal{A}}$  and  $(h_i)_{i \notin \mathcal{A}}$  then, the following equality holds

$$\sum_{i \in \mathcal{A}} S(\beta_i^*) h_i = h'_{\mathcal{A}} S(\beta_{\mathcal{A}}^*) = h'_{\mathcal{A}} X'_{\mathcal{A}} X_{\mathcal{A}} (X'_{\mathcal{A}} X_{\mathcal{A}})^{-1} S(\beta_{\mathcal{A}}^*).$$

Because  $\mathbf{0} = Xh = X_{\mathcal{A}} h_{\mathcal{A}} + X_{\bar{\mathcal{A}}} h_{\bar{\mathcal{A}}}$ , one deduces the following inequalities

$$\begin{aligned} |h'_{\mathcal{A}} S(\beta_{\mathcal{A}}^*)| &= |h'_{\bar{\mathcal{A}}} X'_{\bar{\mathcal{A}}} X_{\mathcal{A}} (X'_{\mathcal{A}} X_{\mathcal{A}})^{-1} S(\beta_{\mathcal{A}}^*)|, \\ &\leq \|h_{\bar{\mathcal{A}}}\|_1 \|X'_{\bar{\mathcal{A}}} X_{\mathcal{A}} (X'_{\mathcal{A}} X_{\mathcal{A}})^{-1} S(\beta_{\mathcal{A}}^*)\|_{\infty}. \end{aligned} \quad (24)$$

Consequently, when the irrerepresentable condition holds for  $\beta^*$  namely, when  $\|X'_{\bar{\mathcal{A}}} X_{\mathcal{A}} (X'_{\mathcal{A}} X_{\mathcal{A}})^{-1} S(\beta_{\mathcal{A}}^*)\|_{\infty} \leq 1$  then, the inequality (24) gives  $|h'_{\mathcal{A}} S(\beta_{\mathcal{A}}^*)| \leq \|h_{\bar{\mathcal{A}}}\|_1$ . Thus, by the equivalence given in (23),  $\beta^*$  is a solution of the following basis pursuit problem

$$\text{minimize } \|\gamma\|_1 \text{ subject to } X\gamma = X\beta^*$$

Because  $X$  is in general position the previous optimisation problem has a unique solution (see *e.g.* the proposition 1 in appendix) thus  $X\beta^* = X\gamma$  and  $\gamma \neq \beta^*$  implies that  $\|\gamma\|_1 > \|\beta^*\|_1$  namely  $\beta^*$  is identifiable with respect to the  $l^1$  norm.  $\square$

## Supplementary material

We already said that when  $X$  is in general position the minimizer of the problem (??) is unique, we also stressed that the estimator derived by minimizing (??) when  $R > 0$  is a LASSO. When the LASSO is written in usual way as in (2), a sketch of proof given in Tibshirani [27] shows the uniqueness of the LASSO estimator when  $X$  is in general position. In order to provide a self content article, we show that when  $X$  is in general position the minimizer of the problem (??) is unique when  $R = 0$  as well as when  $R > 0$ . We already stressed that when  $\beta^*$  is identifiable with respect to the  $l^1$  norm then  $\beta^*$  is sparse. We are going to show that when the identifiability holds for  $\beta^*$  then the family  $(X_i)_{i \in \text{supp}(\beta^*)}$  is linearly independent and thus the number of components of  $\beta^*$  equal to 0 is larger than  $p - n$ .

## References

- [1] R. F. Barber and E. J. Candès. Controlling the false discovery rate via knockoffs. *The Annals of Statistics*, 43(5):2055–2085, 2015.
- [2] Mohsen Bayati and Andrea Montanari. The LASSO risk for Gaussian matrices. *IEEE Transactions on Information Theory*, 58(4):1997–2017, 2012.
- [3] Dimitri P Bertsekas. *Nonlinear programming*. Athena scientific Belmont, 1999.
- [4] M. Bogdan, E. J. Candès, W. Su, and A. Weinstein. Off the beaten path: ranking variables with cross-validated lasso. *Technical Report, University of Wroclaw*, 2018.
- [5] Peter Bühlmann and Sara van de Geer. *Statistics for High-Dimensional Data: Methods, Theory and Applications*. Springer, 2011.
- [6] T Tony Cai and Anru Zhang. Sharp rip bound for sparse signal and low-rank matrix recovery. *Applied and Computational Harmonic Analysis*, 35(1):74–93, 2013.
- [7] Emmanuel J Candès. The restricted isometry property and its implications for compressed sensing. *Comptes Rendus Mathématique*, 346(9):589–592, 2008.
- [8] Emmanuel J. Candès, Yingying Fan, Lucas Janson, and Jinchi Lv. Panning for gold: Model-free knockoffs for high-dimensional controlled variable selection. *arXiv preprint arXiv:1610.02351*, 2016. To appear in Journal of the Royal Statistical Society Series B.
- [9] Emmanuel J Candès, Justin K Romberg, and Terence Tao. Stable signal recovery from incomplete and inaccurate measurements. *Communications on pure and applied mathematics*, 59(8):1207–1223, 2006.

- [10] Scott Shaobing Chen, David L Donoho, and Michael A Saunders. Atomic decomposition by basis pursuit. *SIAM review*, 43(1):129–159, 2001.
- [11] Ingrid Daubechies, Ronald DeVore, Massimo Fornasier, and C Sinan Güntürk. Iteratively reweighted least squares minimization for sparse recovery. *Communications on pure and applied mathematics*, 63(1):1–38, 2010.
- [12] Pascaline Descloux and Sylvain Sardy. Model selection with lasso-zero: adding straw to the haystack to better find needles. *arXiv preprint arXiv:1805.05133*, 2018.
- [13] D. L. Donoho and J. Tanner. Observed universality of phase transitions in high-dimensional geometry, with implications for modern data analysis and signal processing. *Philosophical Trans. R. Soc. A*, 367(1906):4273–4293, 2009.
- [14] David L Donoho and Michael Elad. Optimally sparse representation in general (nonorthogonal) dictionaries via  $\ell_1$  minimization. *Proceedings of the National Academy of Sciences*, 100(5):2197–2202, 2003.
- [15] David L Donoho and Jared Tanner. Precise undersampling theorems. *Proceedings of the IEEE*, 98(6):913–924, 2010.
- [16] Charles Dossal. A necessary and sufficient condition for exact sparse recovery by  $\ell_1$  minimization. *Comptes Rendus Mathématique*, 350(1):117–120, 2012.
- [17] Charles Dossal, Marie-Line Chabanol, Gabriel Peyré, and Jalal Fadili. Sharp support recovery from noisy random measurements by  $\ell_1$ -minimization. *Applied and Computational Harmonic Analysis*, 33(1):24–43, 2012.
- [18] Simon Foucart and Holger Rauhut. *A mathematical introduction to compressive sensing*, volume 1. Springer, 2013.
- [19] Caroline Giacobino, Sylvain Sardy, Jairo Diaz-Rodriguez, Nick Hengartner, et al. Quantile universal threshold. *Electronic Journal of Statistics*, 11(2):4701–4722, 2017.
- [20] Rémi Gribonval and Morten Nielsen. Sparse representations in unions of bases. *IEEE Transactions on Information Theory*, 49(12):3320–3325, 2003.
- [21] Nicolai Meinshausen and Peter Bühlmann. High-dimensional graphs and variable selection with the lasso. *The Annals of Statistics*, 34(3):1436–1462, 2006.
- [22] Hatef Monajemi, Sina Jafarpour, Matan Gavish, David L Donoho, Sivaram Ambikasaran, Sergio Bacallado, Dinesh Bharadia, Yuxin Chen, Young Choi, Mainak Chowdhury, et al. Deterministic matrices matching the compressed sensing phase transitions of gaussian random matrices. *Proceedings of the National Academy of Sciences*, 110(4):1181–1186, 2013.

- [23] Weijie J Su, Małgorzata Bogdan, and Emmanuel J. Candès. False discoveries occur early on the lasso path. *The Annals of Statistics*, 45(5):2133–2150, 2017.
- [24] Patrick Tardivel. *Représentation parcimonieuse et procédures de tests multiples: application à la métabolomique*. PhD thesis, Université de Toulouse, Université Toulouse III-Paul Sabatier, 2017.
- [25] Patrick JC Tardivel, Rémi Servien, and Didier Concordet. Sparsest representations and approximations of an underdetermined linear system. *Inverse Problems*, 34(5):055002, 2018.
- [26] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):267–288, 1996.
- [27] Ryan J Tibshirani et al. The lasso problem and uniqueness. *Electronic Journal of Statistics*, 7:1456–1490, 2013.
- [28] Martin J Wainwright. Sharp thresholds for high-dimensional and noisy sparsity recovery using constrained quadratic programming (lasso). *IEEE transactions on information theory*, 55(5):2183–2202, 2009.
- [29] Asaf Weinstein, Rina Barber, and Emmanuel J. Candès. A power and prediction analysis for knockoffs with lasso statistics. *arXiv preprint arXiv:1712.06465*, 2017.
- [30] Shuaiwen Wang Haolei Weng and Arian Maleki. Which bridge estimator is the best for variable selection? *arxiv*, 2018.
- [31] Peng Zhao and Bin Yu. On model selection consistency of lasso. *The Journal of Machine Learning Research*, 7:2541–2563, 2006.
- [32] Hui Zou. The adaptive lasso and its oracle properties. *Journal of the American statistical association*, 101(476):1418–1429, 2006.