

# Multiple Testing and Variable Selection along Least Angle Regression's path

Jean-Marc Azaïs<sup>•</sup> and Yann De Castro<sup>\*</sup>

<sup>•</sup>*Institut de Mathématiques de Toulouse  
Université Paul Sabatier, 118 route de Narbonne, 31062 Toulouse, France*

<sup>\*</sup>*CERMICS  
École des Ponts ParisTech, 77455 Marne la Vallée, France*

**Abstract:** In this article we investigate the outcomes of the standard Least Angle Regression (LAR) algorithm in high dimensions under the Gaussian noise assumption. We give the exact law of the sequence of knots conditional on the sequence of variables entering the model, i.e., the “post-selection” law of the knots of the LAR. Based on this result, we prove an exact of the False Discovery Rate (FDR) in the orthogonal design case and an exact control of the existence of false negatives in the general design case.

First, we build a sequence of testing procedures on the variables entering the model and we give an exact control of the FDR in the orthogonal design case when the noise level can be unknown. Second, we introduce a new exact testing procedure on the existence of false negatives when the noise level can be unknown. This testing procedure can be deployed after any support selection procedure that will produce an estimation of the support (i.e., the indexes of nonzero coefficients) for any designs. The type *I* error of the test can be exactly controlled as long as the selection procedure follows some elementary hypotheses, referred to as “admissible selection procedures”. These support selection procedures are such that the estimation of the support is given by the  $k$  first variables entering the model where the random variable  $k$  is a stopping time. Monte-Carlo simulations and a real data experiment are provided to illustrate our results.

**MSC 2010 subject classifications:** Primary 62E15, 62F03, 60G15, 62H10, 62H15; secondary 60E05, 60G10, 62J05, 94A08.

**Keywords and phrases:** Multiple Testing, False Discovery Rate, High-Dimensional Statistics, Post-Selection Inference.

*Preprint version of June 27, 2019*

## 1. Introduction

Parsimonious models have become an ubiquitous tool to tackle high-dimensional representations with a small budget of observations. Successful applications is signal processing (see for instance the pioneering works [13, 12] and references therein), biology (see for instance [3] or [11, Chapter 1.4] and references therein)... have shown that the existence of an (almost) sparse representations in some well chosen basis is reasonable assumption in practice. These important successes have put focus on High-Dimensional Statistics in the past decades and they may be due to the deployment of tractable algorithms with strong theoretical guarantees. Among the large panoply of methods, we have seen emerged  $\ell_1$ -regularization which may have found a fine balance between tractability and performances. Nowadays, sparse regression techniques based on  $\ell_1$ -regularization are a common and powerful tool in high-dimensional settings. Popular estimators, among which one may point LASSO [27] or SLOPE [10], are known to achieve minimax rate of prediction and to satisfy sharp oracle inequalities under conditions on the design such as Restricted Eigenvalue [7, 4] or Compatibility [11, 30]. Recent avances have focused on a deeper understanding of these techniques looking at confidence intervals and testing procedures (see [30, Chapter 6] and references therein) or false discovery rate control (e.g., [3]) for instance. These new results aim at describing the (asymptotic or non asymptotic) law of the outcomes of  $\ell_1$ -minimization regression. This line of works addresses important issues encountered in practice. Assessing the uncertainty of popular estimators give strong guarantees on the estimation produced, e.g., the false discovery rate is controlled or a confidence interval on linear statistics of the estimator can be given.

### 1.1. Least Angle Regression algorithm, Support Selection and FDR

Least Angle Regression (LAR) algorithm has been introduced in the seminal article [14]. This forward procedure produces a sequence of knots  $\lambda_1, \lambda_2, \dots$  based on a control of the residuals in  $\ell_\infty$ -norm. This sequence of knots is closely related to the sequence of knots of LASSO [27], as they differ by only one rule: “Only in the LASSO case, if a nonzero coefficient crosses zero before the next variable enters, drop it from the active set and recompute the current joint least-squares direction”, as mentioned in [28, Page 120] for instance. This paper focuses on the LAR algorithm and presents three useful equivalent formulations of the LAR algorithm, see Section A.1. As far as we know, the last formulation given by Algorithm 1, based on a recursive formulation, is new.

One specific task is to estimate the support of the target sparse vector, namely identify the true positives. In particular, one may take the support of a LASSO (or SLOPE) solution as an estimate of the support solution. This strategy has been intensively studied in the literature, one may consider [32, 10, 30, 4] and references therein. Support selection has been studied under the so-called “Irrepresentable Condition” (IC), as presented for instance in the books [30, Page 53] and [11, Sec. 7.5.1] and also referred to as the “Mutual Incoherence Condition” [32]. Under the so-called “Beta-Min Condition”, one may prove [11, 30] that the LASSO asymptotically returns the true support. In this article, we investigate the existence of false non-negatives and we present exact non-asymptotic testing procedures to handle this issue, see Section 2.4. Our procedure does not require IC but a much weaker condition, referred to as the “Empirical Irrepresentable Check” (See Section 2.1.3), that can be checked in polynomial time.

Another recent issue is to control the *False Discovery Rate* (FDR) in high-dimensional setting, as for instance in [3] and references therein; or the *Joint family-wise Error Rate* as in [8] and references therein. In this paper, we investigate the consecutive *spacings* of knots of the LAR as testing statistics and we prove an exact FDR control using a Benjamini–Hochberg procedure [5] in the orthogonal design case, see Section 2.3. Our proof (see Appendix C.7) is based on the *Weak Positive Regression Dependency* (WPRDS), the reader may consult [9] or the survey [23], and Knothe–Roseblatt transport, see for instance [24, Sec.2.3, P.67] or [31, P.20], which is based on conditional quantile transforms.

### 1.2. Post-Selection Inference in High-Dimensions

This paper presents a class of new tests issued from  $\ell_1$ -minimization regression in high-dimensions. Following the original idea of [20], we study tests based on the knots of the LAR’s path. Note that, conditional on the sequence of indexes selected by LAR, the law of two consecutive knots has been studied for the first time by [20] referred to as the *Spacing test* [29]. Later, the article [1] proved that this test is unbiased and introduce a *studentized* version of this test. On the same note, inference after model selection has been studied in several papers, as for instance [15, 25] or [26] for *selective inference* and a joint estimation of the noise level. This line of works studies a single test on a linear statistics while one may ask for a simultaneous control of several tests as in multiple testing frame. To the best of our knowledge, this paper is the first to study the joint law of multiple spacing tests of LAR’s knots in a non-asymptotic frame, see Sections 2.2 and 2.3.

One may point others approaches for building confidence intervals and/or testing procedures in high-dimensional settings as follows. Simultaneous controls of confidence intervals independently of the selection procedure have been studied under the concept of *post-selection constants* as introduced in [6] and studied for instance in [2]. Asymptotic confidence intervals can be build using the *de-sparsified LASSO*, the reader may refer to [30, Chapter 5] and references therein. We also point a recent study [18] of the problem of FDR control as the sample size tends to infinity using *debiased LASSO*.

---

**Algorithm 1:** LAR algorithm (“recursive” formulation)

---

**Data:** Correlations vector  $\bar{Z}$  and variance-covariance matrix  $\bar{R}$ .  
**Result:** Sequence  $((\lambda_k, \bar{v}_k, \varepsilon_k))_{k \geq 1}$  where  $\lambda_1 \geq \lambda_2 \geq \dots > 0$  are the knots, and  $\bar{v}_1, \bar{v}_2, \dots$  are the variables that enter the model with signs  $\varepsilon_1, \varepsilon_2, \dots$  ( $\varepsilon_k = \pm 1$ ).

*/\* Define the recursive function Rec() that would be applied repeatedly. The inputs of Rec() are Z a vector, R a SDP matrix and T a vector. \*/*

**Function** Rec( $R, Z, T$ ):

Compute

$$\lambda = \max_{\{j: T_j < 1\}} \left\{ \frac{Z_j}{1 - T_j} \right\} \text{ and } i = \arg \max_{\{j: T_j < 1\}} \left\{ \frac{Z_j}{1 - T_j} \right\}.$$

*/\* The following recursions are given by (24), (25) and (26). \*/*

Update

$$\begin{aligned} \mathbf{x} &= R_{i\cdot} / R_{ii} \\ R &\leftarrow R - \mathbf{x} R_{i\cdot}^\top \\ Z &\leftarrow Z - \mathbf{x} Z_i \\ T &\leftarrow T + \mathbf{x}(1 - T_i) \end{aligned}$$

**return** ( $R, Z, T, \lambda, i$ )

1 Set  $k = 0, T = 0, Z = (\bar{Z}, -\bar{Z})$  and  $R$  as in (10).  
*/\* Use the following recursion function to compute the LAR path. \*/*  
2 Update  $k \leftarrow k + 1$  and compute  $(R, Z, T, \lambda_k, \hat{v}_k) = \text{Rec}(R, Z, T)$   
Set  $\bar{v}_k = \hat{v}_k \bmod p$  and  $\varepsilon_k = 1 - 2(\hat{v}_k - \bar{v}_k)/p \in \pm 1$ .

---

### 1.3. Outline of the paper

In Section 2, we introduce our method along with the framework and the notion of “*Empirical Irrepresentable Check*” (EIC), a condition that can be checked in polynomial time, see Section 2.1.3. We divided our contributions in three part: Section 2.2 presents the joint law of the knots of LAR under EIC; Section 2.3 gives an exact control of FDR in the *orthogonal design case*; and Section 2.4 presents a general framework to exactly detect false negatives for general designs and some selection procedures referred to as “*admissible*”.

Detailed mathematical statements are given in Section 3 including a method to “studentize” all the tests by an independent estimation of variance.

Section 4 presents practical applications including an illustration on real data.

The Appendix gives details on the different algorithms to compute LAR, proofs of the statements and a list of notation.

## 2. Exact Controls using Least Angle Regression, an introductory presentation

### 2.1. Notation, LAR formulations and Assumptions

#### 2.1.1. Least Angle Regression (LAR)

Consider the linear model in high-dimensions where the number of observations  $n$  may be less than the number of predictors  $p$ . Consider the Least Angle Regression (LAR) algorithm where we denote by  $(\lambda_k)_{k \geq 1}$  the sequence of knots and by  $(\bar{v}_k, \varepsilon_k)_{k \geq 1}$  the sequence of variables  $\bar{v}_k \in [p]$  and signs  $\varepsilon_k \in \{\pm 1\}$  that enter the model along the LAR path, see for instance [28, Chapter 5.6] or [14] for standard description of this algorithm. We recall this algorithm in Algorithm 2 and we present equivalent formulations in Algorithm 3 (using orthogonal projections) and Algorithm 1 (using a recursion). The interested reader may find their analysis in Appendices A and B.

In particular, we present here Algorithm 1 that consists in three lines, applying the same function recursively. We introduce some notation that we will be useful throughout this paper.

We denote by  $(\widehat{v}_1, \dots, \widehat{v}_k) \in [2p]^k$  the “signed” variables that enter the model along the LAR path with the convention that

$$\widehat{v}_\ell := \bar{v}_\ell + p \left( \frac{1 - \varepsilon_\ell}{2} \right), \quad (1)$$

so that  $\widehat{v}_\ell \in [2p]$  is a useful way of encoding both the variable  $\bar{v}_\ell \in [p]$  and its sign  $\varepsilon_\ell = \pm 1$  as used in Algorithm 1. We denote by  $\bar{Z}$  the vector such that  $\bar{Z}_k$  is the scalar product between the  $k$ -th predictor and the response variable, and we denote by  $\bar{R}$  its *correlation* matrix, see (9) and (10) for further details.

### 2.1.2. Nested Models

Our analysis is based conditionally to  $(\widehat{v}_1, \dots, \widehat{v}_n)$  and in this spirit it can be referred to as a “Post-Section” procedure. The selected model  $\widehat{S}$  would be chosen among the family of nested models

$$\underbrace{\{\bar{v}_1\}}_{\bar{S}^1} \subset \underbrace{\{\bar{v}_1, \bar{v}_2\}}_{\bar{S}^2} \subset \underbrace{\{\bar{v}_1, \bar{v}_2, \dots, \bar{v}_k\}}_{\bar{S}^k} \subset \dots \subset \underbrace{\{\bar{v}_1, \bar{v}_2, \dots, \bar{v}_n\}}_{\bar{S}^n}. \quad (2)$$

Respectively, denote

$$E_1 \subset \dots \subset E_k := \text{Span}(X_{\bar{v}_1}, \dots, X_{\bar{v}_k}) \subset \dots \subset E_n = \mathbb{R}^n,$$

the corresponding family of nested subspaces of  $\mathbb{R}^n$ . In the sequel, denote by  $P_k(Y)$  (resp.  $P_k^\perp(Y)$ ) the orthogonal projection of the observation  $Y$  onto  $E_k$  (resp. the orthogonal of  $E_k$ ) for all  $k \geq 1$ .

### 2.1.3. Empirical Irrepresentable Check

In the sequel we will build our testing statistics on the joint law of the first  $K$  knots of the LAR. A detailed discussion on the choice of  $K$  is given in Section 2.5, in particular  $K$  will be such that the so-called “*Empirical Irrepresentable Check of order  $K$* ” holds.

This property will be defined and studied in this section. Note that the “*Irrepresentable Condition*” is a standard condition, as presented for instance in the books [30, Page 53] and [11, Sec. 7.5.1] and also referred to as the “*Mutual Incoherence Condition*” [32].

**Definition 1** (Irrepresentable Condition of order  $K$ ). *The design matrix  $X$  satisfies the Irrepresentable Condition of order  $K$  if and only if*

$$\forall S \subset [p] \text{ s.t. } \#S \leq K, \quad \max_{j \in [p] \setminus S} \max_{\|v\|_\infty \leq 1} X_j^\top X_S (X_S^\top X_S)^{-1} v < 1, \quad (\mathcal{A}_{\text{Irr.}})$$

where  $X_j$  denotes the  $j^{\text{th}}$  column of  $X$  and  $X_S$  the submatrix of  $X$  obtained by keeping the columns indexed by  $S$ .

**Remark 1.** *This condition has been intensively studied in the literature and it is now well established that some random matrix models satisfies it with high probability. For instance, one may refer to the article [32] where it is shown that a design matrix  $X \in \mathbb{R}^{n \times p}$  whose rows are drawn independently with respect to a centered Gaussian distribution with variance-covariance matrix satisfying  $(\mathcal{A}_{\text{Irr.}})$  (for instance the Identity matrix) satisfies  $(\mathcal{A}_{\text{Irr.}})$  with high probability when  $n \gtrsim K \log(p - K)$ , where  $\gtrsim$  denotes an inequality up to some multiplicative constant.*

In practice, the Irrepresentable Condition  $(\mathcal{A}_{\text{Irr.}})$  is a strong requirement on the design  $X$  and, additionally, this condition cannot be checked in polynomial time. One important feature of our result, see Theorem 3, is that we do not require Irrepresentable Condition but a much weaker requirement referred to as the “*Empirical Irrepresentable Check*” of order  $K$ .

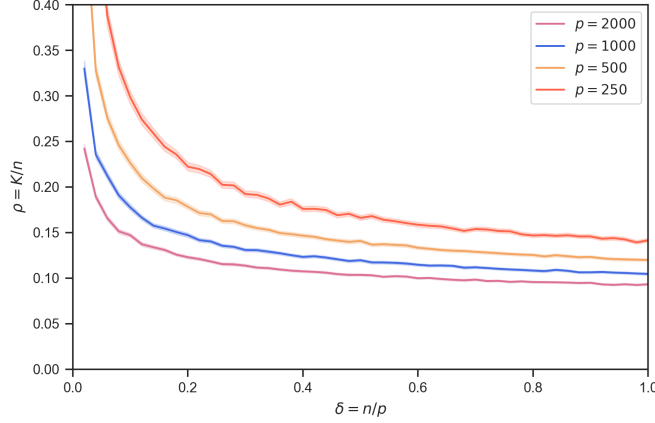


Figure 1. Observed empirical “phase transition” over 1,000 Monte-Carlo repetitions on the Empirical Irrepresentable Check of order  $K$  with  $\rho := K/n$  and  $\delta := n/p$  for different values of  $p$ . We considered a design  $X \in \mathbb{R}^{n \times p}$  with independent column vectors uniformly distributed on the sphere and an independent  $y \in \mathbb{R}^n$  with i.i.d. standard Gaussian entries, and we computed the variables  $(\hat{v}_1, \dots, \hat{v}_n)$  entering the model with LAR. The plot represents the value  $K$  defined as the largest order for which  $(\hat{\mathcal{A}}_{\text{Irr.}})$  holds with respect to  $(\hat{v}_1, \dots, \hat{v}_n)$ .

**Definition 2** (Empirical Irrepresentable Check of order  $K$  w.r.t  $(\hat{v}_1, \dots, \hat{v}_K)$ ). The design matrix  $X$  satisfies the Empirical Irrepresentable Check of order  $K$  with respect to  $(\hat{v}_1, \dots, \hat{v}_K)$  if and only if for all the supports  $(\bar{S}^k)_{k \in [K]}$  and all the signs  $(\varepsilon^k)_{k \in [K]}$  which enter the model at steps  $k \in [K]$

$$\forall k \in [K], \forall j \notin \bar{S}^k := \{\bar{v}_1, \dots, \bar{v}_k\}, \quad X_j^\top X_{\bar{S}^k} (X_{\bar{S}^k}^\top X_{\bar{S}^k})^{-1} \varepsilon^k < 1 \quad (\hat{\mathcal{A}}_{\text{Irr.}})$$

where  $\bar{S}^k$  is the selected support at step  $k$  (see (2)) and  $\varepsilon^k := \{\varepsilon_1, \dots, \varepsilon_k\}$  are the corresponding signs of the selected variables.

**Remark 2.** Observe that Irrepresentable Condition  $(\mathcal{A}_{\text{Irr.}})$  of order  $K$  implies Empirical Irrepresentable Check  $(\hat{\mathcal{A}}_{\text{Irr.}})$  of order  $K$  with respect to any possible  $(\hat{v}_1, \dots, \hat{v}_K)$ .

**Remark 3.** Note that this condition can be checked in  $\mathcal{O}(n \times p \times K + n \times K^{c_{X, \text{inv}}})$  time where  $\mathcal{O}(n \times K^{c_{X, \text{inv}}})$  accounts for the computational time of the  $K$  operations  $X_S^\top (X_S^\top X_S)^{-1} e$  where  $e \in \mathbb{R}^k$  is a fixed vector and  $S \subseteq [p]$  has size  $k$  for  $k = 1, \dots, K$ . A standard bound for this last term would be  $c_{X, \text{inv}} \leq 6.3728639$ , see for instance [19].

**Remark 4.** When computing the LAR path, one has to compute the values  $X_j^\top X_{\bar{S}^k} (X_{\bar{S}^k}^\top X_{\bar{S}^k})^{-1} \varepsilon^k$ , see for instance Algorithm 2 or Algorithm 3 where these values are given by  $\theta$  as shown by Proposition 2. It means that, in practice, one has for free the maximal order  $K$  for which Empirical Irrepresentable Check holds. Furthermore, this latter remark shows that the maximal order  $K$  for which Empirical Irrepresentable Check holds is a statistic of the variables  $\bar{v}_k \in [p]$  and signs  $\varepsilon_k \in \{\pm 1\}$  that enter the model along the LAR path.

In Figure 1, in the particular case where the design matrix  $X \in \mathbb{R}^{n \times p}$  is constructed from independent column vectors uniformly distributed on the sphere and an independent  $y \in \mathbb{R}^n$  with i.i.d. standard Gaussian entries, and we computed the maximal order  $K$  for which the Empirical Irrepresentable Check holds. Due to symmetry of the law of the design and of the independent observation  $y$ , note that the probability to observe a given path  $\mathbb{P}\{\hat{v}_1 = v_1, \dots, \hat{v}_n = v_n\}$  is almost the same for all possible paths, namely  $\simeq 1/(n!)$ . Then we can view the plot in Figure 1 as the largest value  $K$  for which  $(\hat{\mathcal{A}}_{\text{Irr.}})$  holds when the support  $\{\hat{v}_1 = v_1, \dots, \hat{v}_K = v_K\}$  is uniformly distributed.

In Figure 1, we chose  $p = 250, 500, 1000, 2000$ , different values of  $n$  and we plot a shaded blue line representing 95% of the values  $\rho$  (among 1,000 repetitions). Interestingly, these values are

concentrated around the solid blue line referred to as the “empirical phase transition” for the Gaussian model.

**Example 5.** For example, we found that for  $p = 1,000$  and  $n = 300$  ( $\delta = n/p = 0.3$ ), Empirical Irrepresentable Check ( $\hat{\mathcal{A}}_{\text{Irr.}}$ ) of order  $K$  holds when  $K$  is about  $K \simeq \rho \times 300 = 0.13 \times 300 = 39$ .

For every couple ( $\delta = n/p, \rho = K/n$ ) below this curve we observe that we have ( $\hat{\mathcal{A}}_{\text{Irr.}}$ ) of order  $K$  with overwhelming probability as  $p$  is sufficiently large. We did not pursue this issue in the present paper but we would like to emphasize that ( $\hat{\mathcal{A}}_{\text{Irr.}}$ ) holds empirically for values  $K$  of the order of 10% of  $n$  at least. In addition, on real data of Section 4.2, with about  $p = 200$  variables and  $n = 700$  observations with found  $K$  of the order of 30.

## 2.2. Result one: Mixture of Gaussian Order Statistics in LAR

Recall that this article concerns the linear model in high-dimensions where the number of predictors  $p$  can be much greater than the number of samples  $n$ . We denote by  $Y \in \mathbb{R}^n$  the response variable and we assume that

$$Y = X^0 \beta^0 + \eta \sim \mathcal{N}_n(X^0 \beta^0, \sigma^2 \Sigma), \quad (3)$$

where  $\eta \sim \mathcal{N}_n(0, \sigma^2 \Sigma)$  is some Gaussian noise, the matrix  $\Sigma \succ 0$  is some known positive definite matrix, the noise level  $\sigma > 0$  may be known or that has to be estimated depending on the context, and  $X^0 \in \mathbb{R}^{n \times p}$ . In this paper, we are interested in selecting the true support  $S^0$  of  $\beta^0$ , where the support is defined by

$$S^0 := \{k \in [p] : \beta_k^0 \neq 0\}.$$

Provided that  $\Sigma$  is known, one can always consider the homoscedastic-independent version of the responses, namely

$$\Sigma^{-\frac{1}{2}} Y = X \beta^0 + \Sigma^{-\frac{1}{2}} \eta \sim \mathcal{N}_n(X \beta^0, \sigma^2 \text{Id}_n).$$

where the design  $X$  is given by  $X := \Sigma^{-\frac{1}{2}} X^0$ . From the observation  $\Sigma^{-\frac{1}{2}} Y$  and the design  $X$  (or equivalently from  $(\bar{Z}, \bar{R})$  defined in (9)), one can compute the Least Angle Regression (LAR) knots  $(\lambda_k)_{k \geq 1}$  and the sequence  $(\bar{i}_k, \varepsilon_k)_{k \geq 1}$  of variables  $\bar{i}_k \in [p]$  and signs  $\varepsilon_k \in \{\pm 1\}$ . In the sequel, the law (resp. the conditional law) of a random vector  $V$  (resp. conditionally to  $W$ ) is denoted by  $\mathcal{L}(V)$  (resp.  $\mathcal{L}(V|W)$ ). One of the main discovery of this paper (see Theorem 3) is the following:

*Under Empirical Irrepresentable Check ( $\hat{\mathcal{A}}_{\text{Irr.}}$ ) of order  $K$ , the LAR knots  $(\lambda_1, \dots, \lambda_K)$  behave as a mixture of Gaussian order statistics with heterogeneous variances, namely it holds*

$$\begin{aligned} & \mathcal{L}(\lambda_1, \dots, \lambda_K | \lambda_{K+1}) \\ &= \sum_{(i_1, \dots, i_K) \in [2p]^K} \pi_{(i_1, \dots, i_K, \lambda_{K+1})} \underbrace{Z_{(i_1, \dots, i_K, \lambda_{K+1})}^{-1} \left[ \bigotimes_{k=1}^K \gamma_{m_k, v_k^2} \right] \mathbb{1}_{\{\lambda_1 \geq \dots \geq \lambda_K \geq \lambda_{K+1}\}}}_{\mathcal{L}(\lambda_1, \dots, \lambda_K | \hat{i}_1 = i_1, \dots, \hat{i}_K = i_K, \lambda_{K+1})} \end{aligned} \quad (4)$$

where  $\pi_{(i_1, \dots, i_K, \lambda_{K+1})} = \mathbb{P}\{\hat{i}_1 = i_1, \dots, \hat{i}_K = i_K | \lambda_{K+1}\}$  can be interpreted as mixing probabilities,  $Z_{(i_1, \dots, i_K, \lambda_{K+1})}$  is a normalizing constant, and  $\gamma_{m_k, v_k^2}$  is a short notation for the Gaussian law such that

- its mean  $m_k$  is given by (15) and depends only on the covariance  $\sigma^2 \Sigma$ , the selected support  $(\bar{i}_1, \dots, \bar{i}_k)$  at step  $k$ , and the orthogonal projection of  $\bar{\mu}^0 := X^\top X \beta^0$  given by (5);
- and its variance  $v_k^2 := \sigma^2 \rho_k^2$  is given by (13) and depends only on the covariance  $\sigma^2 \Sigma$  and the selected support  $(i_1, \dots, i_k)$  at step  $k$ .

As presented in the sequel, this result is the key step to prove an exact control of the FDR, see Section 2.3, and an exact control of the False Negatives of some selection procedures, see Section 2.4. These controls are obtained from  $p$ -values  $\hat{\alpha}_{a,b,c}$  described in Theorem 5 and (19). We will see that **the  $p$ -value  $\hat{\alpha}_{abc}$  detects abnormal large values of  $\lambda_b$  conditional on  $(\lambda_a, \lambda_c)$**  and that, in the orthogonal design case, Theorem 6 shows that the test based on  $\hat{\alpha}_{a,a+1,K+1}$  is uniformly more powerful than tests based on  $\hat{\alpha}_{x,y,z}$  with  $a \leq x < y < z \leq K+1$ .

### 2.3. Result two: Control of False Discovery Rate in the Orthogonal Design case

#### 2.3.1. Presentation in the general case

Assume that Empirical Irrepresentable Check ( $\hat{\mathcal{A}}_{\text{irr.}}$ ) of order  $K$  holds. The sequence of testing procedures under consideration is given by the sequence of  $p$ -values  $(\hat{\alpha}_{k-1,k,k+1})_{k \in [K]}$  that are referred to as “*Spacing tests*” [29] in the literature, they are recalled in Remark 10 below. These  $p$ -values account for **“abnormally large” values of  $\lambda_k$  conditional on  $(\lambda_{k-1}, \lambda_{k+1})$** . Invoking (4), this conditional law of  $\lambda_k$  depends on two unobserved values: the mean  $m_k$ , given by (15), and the variance  $\sigma^2$  (in the case where the variance may be unknown). As discussed in Section 2.4.1, the variance  $\sigma^2$  will be estimated on residuals as described by the heuristic of Remark 6. It entails that one may efficiently estimate the variance and, as we will see in Section 3.4, one can plug this independent estimator of the variance so as to get a studentized version of the testing procedures described here. For sake of readability, we will assume that  $\sigma$  is known and we refer to Section 3.4 for details on how to studentize our procedure. We understand that the law of testing statistics are parametrized by the hypotheses  $(m_k)_{k \in [K]}$ , where  $m_k$  is given by (15).

We recall that we denote  $\bar{\mu}^0 = X^\top X \beta^0$  and  $\bar{\mu}_i^0$  its  $i$ th coordinate. Assuming that predictors are normalised, in the general case, this quantity is the sum of  $\beta_j^0$ 's whose predictors  $X_j$  are highly correlated with the predictor  $X_i$ . Now, given  $\bar{i}_1, \dots, \bar{i}_k \in [p]$  and signs  $\varepsilon_1, \dots, \varepsilon_k \in \{\pm 1\}^k$ , we denote by  $(P_{\bar{i}_1, \dots, \bar{i}_{k-1}}^\perp(\bar{\mu}^0))_{\bar{i}_k}$  the orthogonal projection given by

$$(P_{\bar{i}_1, \dots, \bar{i}_{k-1}}^\perp(\bar{\mu}^0))_{\bar{i}_k} := \varepsilon_k X_{\bar{i}_k}^\top \left[ \text{Id}_n - X_{\bar{S}^{k-1}} \text{diag}(\varepsilon_{\bar{S}^{k-1}}) (X_{\bar{S}^{k-1}}^\top X_{\bar{S}^{k-1}})^{-1} \text{diag}(\varepsilon_{\bar{S}^{k-1}}) X_{\bar{S}^{k-1}}^\top \right] X \beta^0. \quad (5)$$

where  $\text{diag}(\varepsilon_{\bar{S}^{k-1}}) := \text{diag}(\varepsilon_1, \dots, \varepsilon_{k-1})$ .

The tested hypotheses are conditional on the sequence of variables  $(\bar{i}_1, \dots, \bar{i}_K) \in [p]^K$  and signs  $\varepsilon_1, \dots, \varepsilon_K \in \{\pm 1\}^K$  entering the model. The  $p$ -values under consideration here are given by

- $\hat{p}_1 := \hat{\alpha}_{0,1,2}$  is the  $p$ -value testing  $\mathbb{H}_{0,1}$  : “ $m_1 = 0$ ” namely  $\bar{\mu}_{\bar{i}_1}^0 = 0$ ;
- $\hat{p}_2 := \hat{\alpha}_{1,2,3}$  is the  $p$ -value testing  $\mathbb{H}_{0,1}$  : “ $m_2 = 0$ ” namely  $(P_1^\perp(\bar{\mu}^0))_{\bar{i}_2} = 0$ ;
- $\hat{p}_3 := \hat{\alpha}_{2,3,4}$  is the  $p$ -value testing  $\mathbb{H}_{0,1}$  : “ $m_1 = 0$ ” namely  $(P_2^\perp(\bar{\mu}^0))_{\bar{i}_3} = 0$ ;
- and so on...

We write  $I_0$  for the set

$$I_0 = \{k \in [K] : \mathbb{H}_{0,k} \text{ is true}\},$$

Given a subset  $\hat{R} \subseteq [K]$  of hypotheses that we consider as rejected, we call *false positive* (FP) and *true positive* (TP) the quantities  $\text{FP} = \text{card}(\hat{R} \cap I_0)$  and  $\text{TP} = \text{card}(\hat{R} \setminus I_0)$ .

Denote by  $\hat{p}_{(1)} \leq \dots \leq \hat{p}_{(K)}$  the  $p$ -values ranked in a nondecreasing order. Let  $\alpha \in (0, 1)$  and consider the Benjamini-Hochberg procedure, see for instance [5], defined by a rejection set  $\hat{R} \subseteq [K]$  such that  $\hat{R} = \emptyset$  when  $\{k \in [K] : \hat{p}_{(k)} \leq \alpha k/K\} = \emptyset$  and

$$\hat{R} = \{k \in [K] : \hat{p}_k \leq \alpha \hat{k}/K\} \quad \text{where} \quad \hat{k} = \max \{k \in [K] : \hat{p}_{(k)} \leq \alpha k/K\}. \quad (7)$$



Recall the definition of FDR as the mean of False Discovery Proportion (FDP), namely

$$\text{FDR} := \mathbb{E} \left[ \underbrace{\frac{\text{FP}}{\text{FP} + \text{TP}} \mathbf{1}_{\text{FP} + \text{TP} \geq 1}}_{\text{FDP}} \right],$$

where the expectation is unconditional to the sequence of variables entering the model, while the hypotheses that are being tested are conditional on the sequence of variables entering the model. This FDR can be understood invoking the following decomposition

$$\text{FDR} = \sum_{(i_1, \dots, i_K) \in [p]^K} \bar{\pi}_{(i_1, \dots, i_K)} \mathbb{E}[\text{FDP} | \bar{i}_1 = i_1, \dots, \bar{i}_K = i_K],$$

where  $\bar{\pi}_{(i_1, \dots, i_K)} = \mathbb{P}\{\bar{i}_1 = i_1, \dots, \bar{i}_K = i_K\}$ .

### 2.3.2. FDR control of Benjamini–Hochberg procedure in the orthogonal design case

We now consider the orthogonal design case where  $X^\top X = \text{Id}_p$  and the set of  $p$ -values given by (6). Note that  $I_0$  is simply the set of null coordinates of  $\beta$ . Remark also that, Irrepresentable Condition ( $\mathcal{A}_{\text{irr.}}$ ) of order  $p$  holds and so does Empirical Irrepresentable Check ( $\hat{\mathcal{A}}_{\text{irr.}}$ ), see Proposition 2. Note also that  $I_0$  is simply the set of null coordinates of  $\beta$ . It implies that one can consider any value  $K \in [p]$  in the following result.

**Theorem 1.** *Assume that the design is orthogonal, namely it holds  $X^\top X = \text{Id}_p$ , and let  $K \in [p]$ . Let  $(\bar{i}_1, \dots, \bar{i}_K)$  be the first variables entering along the LAR’s path. Consider the  $p$ -values given by (6) and the set  $\hat{R}$  given by (7). Then*

$$\mathbb{E}[\text{FDP} | \bar{i}_1 = i_1, \dots, \bar{i}_K = i_K] \leq \alpha,$$

and so FDR is upper bounded by  $\alpha$ .

The proof of this result is given in Appendix C.7.

One interpretation of *post-selection type* may be given as follows: if one looks at all the experiments giving the same sequence of variables entering the model  $\{\bar{i}_1 = i_1, \dots, \bar{i}_K = i_K\}$  and if one considers the Benjamini–Hochberg procedure for the hypotheses described in Section 2.3.1, then the FDR is exactly controlled by  $\alpha$ .

## 2.4. Result three: Exact Testing Procedure on False Negatives

From the result of Section 2.2, one can present a method to select a model and propose an exact test on false negatives in the general design case. More precisely, we assume here for simplicity that the design  $X$  has full row rank  $n$ . In theory we are able to compute the LAR knots up to  $\lambda_n$ . Of course this is often out of reach due to numerical issues and one has to stop earlier to avoid bad conditioning, but we forget this possibility to avoid heavy notation. One can adapt easily the following text to take into account this possibility. Let us compute the LAR knots  $(\lambda_1, \dots, \lambda_n)$  and the “signed” variables  $(\hat{i}_1, \dots, \hat{i}_n) \in [2p]^n$  entering the model as in (1).

### 2.4.1. Admissible Support Selection Procedures

In order to select  $\hat{S}$ , one may be interested in an estimation  $\hat{\sigma}_{\text{select}}$  of the noise level since, in practice, one usually does not know it. Our procedure is valid as long as the following property ( $\mathcal{P}_1$ ) is satisfied

**Noise Level Estimation for Selection ( $\mathcal{P}_1$ ):** *The estimated noise level  $\hat{\sigma}_{\text{select}}$  is a measurable function of  $P_{n_{\text{select}}}^\perp(Y)$  where  $n_{\text{select}}$  may depend only on  $(\hat{i}_1, \dots, \hat{i}_n)$ .*



In practice, one can chose a fixed  $n_{\text{select}} = n - B$  giving a fixed budget of  $B \geq 1$  independent observations to estimate the variance<sup>1</sup>. Obviously, if the noise level is known, one can simply take  $n_{\text{select}} = n$  and  $\hat{\sigma}_{\text{select}} = \sigma$ .

**Remark 6.** Observe that if the true support  $S^0$  is included in  $\{\bar{i}_1, \bar{i}_2, \dots, \bar{i}_k\}$  then  $P_k^\perp(Y)$  is centered Gaussian vector with known covariance (up to  $\sigma^2$ ) and it is standard to compute an estimator of the noise level  $\sigma$ . In particular, the information contained in  $P_k^\perp(Y)$  does not depend on the target  $\beta^0$ , since  $P_k^\perp(X^0\beta^0) = 0$ , and is purely due to the “noise” part  $\eta$ , as defined in (3).

Note that, for large enough  $n_{\text{select}}$ , the heuristic described in Remark 6 may be true and it entails that one may efficiently estimate the variance on  $P_{n_{\text{select}}}^\perp(Y)$ .

Now, note that choosing a model  $\hat{S}$  is equivalent to choosing a model size  $\hat{m} \in [n]$  so that

$$\hat{S} = \{\bar{i}_1, \bar{i}_2, \dots, \bar{i}_{\hat{m}}\}.$$

Our procedure is flexible on this point and it allows any choice of  $\hat{m}$  as long as the following property ( $\mathcal{P}_2$ ) is satisfied

**Stopping Rule ( $\mathcal{P}_2$ ):** The estimated model size  $\hat{m}$  is a “stopping time”, i.e.  $\mathbb{1}_{\{\hat{m} \leq k\}}$  is a measurable function of  $(P_k(Y), \hat{\sigma}_{\text{select}})$  for all  $k \in [n_{\text{select}}]$ .

In other words, the decision to select a model of size  $\hat{m} = k$  depends only on the part of the observation  $Y$  explained by the predictors  $X_{\bar{i}_1}, \dots, X_{\bar{i}_k}$ . In the sequel, we will present some examples of such “stopping time” procedures, see Section 3.5. We would like to emphasize that our analysis reveal the following conditional independence

$$\mathcal{L}((P_{n_{\text{select}}}(Y), \hat{\sigma}_{\text{select}})|(\hat{i}_1, \dots, \hat{i}_n)) = \mathcal{L}(P_{n_{\text{select}}}(Y)|(\hat{i}_1, \dots, \hat{i}_n)) \otimes \mathcal{L}(\hat{\sigma}_{\text{select}}|(\hat{i}_1, \dots, \hat{i}_n)) \quad (8)$$

which can be advantageously invoked to build Student-type testing statistics to define  $\hat{m}$ .

#### 2.4.2. Exact Testing Procedure on False Negatives

Once one has selected a model of size  $\hat{m}$ , one may be willing to test if  $\hat{S}$  contains the true support  $S^0$  by considering the null hypothesis  $\mathbb{H}_0 : “S^0 \subseteq \hat{S}”$ , namely there is no false negatives. Equivalently, one aim at testing the null hypothesis

$$\mathbb{H}_0 : “X\beta^0 \in E_{\hat{m}}”,$$

at an exact significance level  $\alpha \in (0, 1)$

In this article, we introduce a new **exact** testing procedure that can be deployed when ( $\mathcal{P}_1$ ) and ( $\mathcal{P}_2$ ) hold, namely an “admissible” selection procedure is used to build  $\hat{S}$ . The variance estimator  $\hat{\sigma}_{\text{select}}$  may have been used in the definition of  $\hat{m}$  and our method forbid to use it again to “studentized” our testing procedure. We need to estimate the variance again to preserve some “independence” in the spirit of (8). As in the previous variance estimation, we advocate a “budget” to this task, namely we use the “strata”  $(n_{\text{test}}, n_{\text{select}}]$  (as in the vocabulary of analysis of variance) to build  $\hat{\sigma}_{\text{test}}^2$  an estimation of the variance. More precisely, we use the orthogonal projection  $P_{(n_{\text{test}}, n_{\text{select}}]}(Y)$  which is the the orthogonal projector onto the space  $F$  defined by

$$E_{n_{\text{select}}} = E_{n_{\text{test}}} \overset{\perp}{\oplus} F.$$

Note that this space is generated by  $X_{\bar{i}_{n_{\text{test}}+1}}, \dots, X_{\bar{i}_{n_{\text{select}}}}$ .

The resulting test is based on  $(\lambda_{\hat{m}}/\hat{\sigma}_{\text{test}}, \dots, \lambda_{n_{\text{test}}}/\hat{\sigma}_{\text{test}})$ . More precisely, we know the law of  $\lambda_{\hat{m}_{\text{select}}+1}/\hat{\sigma}_{\text{test}}$  conditional on  $\hat{m}, \lambda_{\hat{m}}, \lambda_{n_{\text{test}}}$  under  $\mathbb{H}_0$ .

<sup>1</sup> One can also look at the conditioning of the variance-covariance matrix of  $P_k^\perp(Y)$  for  $k = n, n-1, \dots$  to choose  $n_{\text{select}}$ . A more detailed discussion on this issue can be found in Section 2.5.

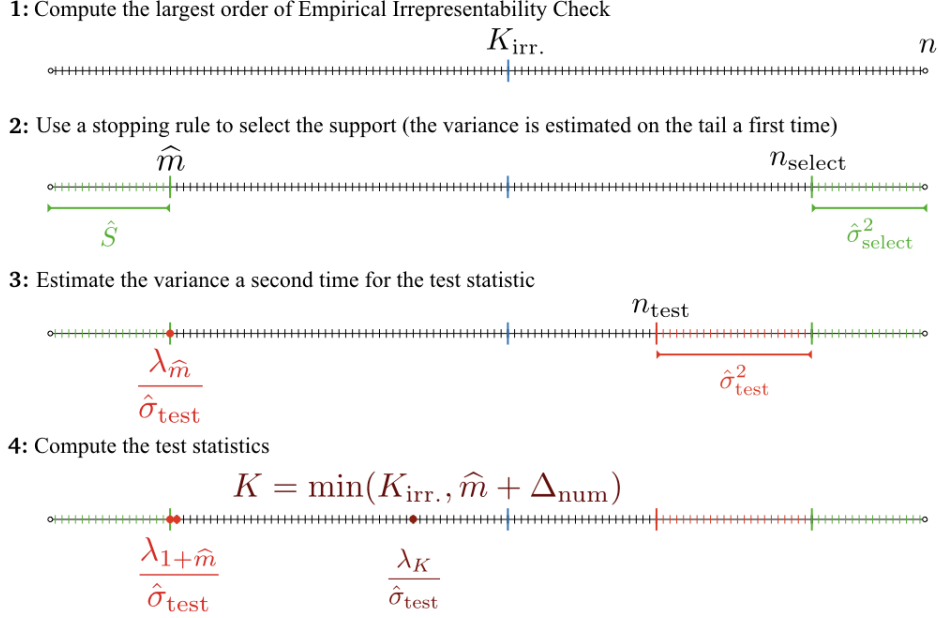


Figure 2. The testing procedure computes the significance of knot  $\lambda_{1+\hat{m}}$  w.r.t. knots  $\lambda_{\hat{m}}$  and  $\lambda_K$  (see Step 4 above), where  $K$  is a compromise between the numerical limitations  $\hat{m} + \Delta_{\text{num}}$  of the cost of integrating a function on  $[0, 1]^{\Delta_{\text{num}}-2}$  and the statistical limitations of the design as depicted by the Empirical Irrepresentable Check.

## 2.5. Practical Framework

We start by presenting the meta-algorithm in which our testing procedure can be deployed. First, we compute the sequence of LAR knots  $(\lambda_1, \dots, \lambda_n)$  and “signed” variables  $(\hat{v}_1, \dots, \hat{v}_n) \in [2p]^n$  entering the model as in (1). Conditionally to the sequence  $\hat{v}_1, \dots, \hat{v}_n$  the quantities displayed in Figure 2 are introduced in the following order.

- Compute  $K_{\text{irr.}}$  which is the maximal value such that  $(\hat{\mathcal{A}}_{\text{irr.}})$  holds true.
- Fixing a budget  $B$  for variance estimation  $\hat{\sigma}_{\text{select}}$  we fix  $n_{\text{select}} = n - B$ . In most of the applications,  $n$  is large and consequently  $n_{\text{select}}$  is also large.
- From the estimation  $\hat{\sigma}_{\text{select}}$  we deduce  $\hat{m}_{\text{select}}$ . Here, the practitioner is free to use any selection procedures as soon as it satisfies  $(\mathcal{P}_1)$  and  $(\mathcal{P}_2)$ .
- Multivariate integration programs have some dimension limitation, this implies that we can only handle simultaneously  $\Delta_{\text{num}}$  knots. As a consequence  $K$  must be smaller or equal than  $\hat{m}_{\text{select}} + \Delta_{\text{num}}$ .
- In most cases the quantity

$$K = \min(K_{\text{irr.}}, \hat{m}_{\text{select}} + \Delta_{\text{num}})$$

is convenient because  $K < n_{\text{select}}$  and there is enough budget in the strata  $(K, n_{\text{select}}]$  to build the second estimator of variance, namely  $\hat{\sigma}_{\text{test}}$ .

- Otherwise we have to make a “trade-off”: diminishing  $K$  to obtain a sufficient budget of degree of freedom in the strata  $(K, n_{\text{select}}]$  to estimate the variance while maintaining  $K$  large to gain power.

An example of values of these parameters is

$$n = 200, \quad p = 300, \quad K_{\text{irr.}} = 35 \text{ (for a Gaussian design), and } \Delta_{\text{num}} = 10.$$

This framework has been deployed on real data in Section 4. Examples of *admissible* procedures satisfying  $(\mathcal{P}_1)$  and  $(\mathcal{P}_2)$  are presented in Section 3.5.

### 3. Main Results

Consider the correlations vector  $\bar{Z}$  of the homoscedastic-independent observation  $\Sigma^{-\frac{1}{2}}Y$  with the design matrix  $X$ , so that

$$\begin{aligned} \bar{Z} &:= X^{0\top} \Sigma^{-1} Y = \bar{R} \beta^0 + X^{0\top} \Sigma^{-1} \eta \sim \mathcal{N}_p(\bar{R} \beta^0, \sigma^2 \bar{R}) \\ \text{and } \bar{R} &:= X^\top X = X^{0\top} \Sigma^{-1} X^0. \end{aligned} \quad (9)$$

From  $\bar{Z}$  and  $\bar{R}$  we can compute the Least Angle Regression (LAR) knots  $(\lambda_k)_{k \geq 1}$  and the sequence  $(\bar{i}_k, \varepsilon_k)_{k \geq 1}$  of variables  $\bar{i}_k$  and signs  $\varepsilon_k$  that enter the model along the LAR path. We will use these statistics as testing statistics. For sake of presentation, we may consider the  $2p$ -vector  $Z := (\bar{Z}, -\bar{Z})$  whose mean is given by  $\mu^0 := (\bar{R} \beta^0, -\bar{R} \beta^0) = (X^\top X \beta^0, -X^\top X \beta^0) = (\bar{\mu}^0, -\bar{\mu}^0)$  and its variance-covariance matrix is  $\sigma^2 R$  with

$$R = \begin{bmatrix} \bar{R} & -\bar{R} \\ -\bar{R} & \bar{R} \end{bmatrix} = \begin{bmatrix} X^\top X & -X^\top X \\ -X^\top X & X^\top X \end{bmatrix}. \quad (10)$$

#### 3.1. Problem reformulation and Assumptions

Remark that  $Z = (Z_i)_i$  is a Gaussian vector of size  $2p$  such that  $Z_i = -Z_{i+p}$  where the indices are considered modulo  $2p$ . Note also that we know the covariance of  $Z$  up to some multiplicative constant  $\sigma^2$ , namely the variance-covariance matrix of  $Z$  is given by  $\sigma^2 R$  where  $R \succeq 0$  is known and  $\sigma^2$  is some parameter that may be known or not.

Now, we can define

$$\forall (i_1, \dots, i_k) \in [2p]^k, \quad \theta_j(i_1, \dots, i_k) := (R_{j,i_1} \cdots R_{j,i_k}) M_{i_1, \dots, i_k}^{-1} (1, \dots, 1), \quad (11)$$

where  $(1, \dots, 1)$  is the column vector of size  $k$  whose entries are equal to one and  $\sigma^2 M_{i_1, \dots, i_k}$  is the variance-covariance matrix of the vector  $(Z_{i_1}, \dots, Z_{i_k})$ , and  $(R_{j,i_1} \cdots R_{j,i_k})$  is a row vector of size  $k$ . Note that  $M_{i_1, \dots, i_k}$  is the submatrix of  $R$  obtained by keeping the columns and the rows indexed by  $\{i_1, \dots, i_k\}$ . Remark that

$$\theta_j(i_1, \dots, i_k) = \mathbb{E}(Z_j \mid Z_{i_1} = 1, \dots, Z_{i_k} = 1),$$

when  $\mathbb{E}Z = 0$ . In our context, the Irrepresentable Condition of order  $K$  is given by

$$\forall k \leq K, \forall (i_1, \dots, i_k) \in [2p]^k, \forall j \notin \{i_1, \dots, i_k\}, \quad \theta_j(i_1, \dots, i_k) < 1, \quad (12)$$

where  $\theta_j(i_1, \dots, i_k)$  is given by (11). In Proposition 2 we show that (12) is equivalent to the Irrepresentable Condition  $(\mathcal{A}_{\text{Irr.}})$  on the design matrix  $X$ . A proof can be found in Appendix C.2.

**Proposition 2.** *Assume  $X$  and  $R$  satisfy (10) then the following assumptions are equivalent:*

- the design matrix  $X$  satisfies  $(\mathcal{A}_{\text{Irr.}})$  of order  $K$ ,
- the variance-covariance matrix  $R$  satisfies (12) of order  $K$ .

Furthermore, if one of these two assumptions hold then “Empirical Irrepresentable Check” of order  $K$  holds, namely

$$\max \left[ \max_{j \neq \hat{i}_1} \theta_j(\hat{i}_1), \dots, \max_{j \neq \hat{i}_1, \dots, \hat{i}_k} \theta_j(\hat{i}_1, \dots, \hat{i}_k) \right] < 1$$

which is an equivalent formulation of  $(\hat{\mathcal{A}}_{\text{Irr.}})$  of Definition 2.

**Remark 7.** One may require that the design is “normalized” so that  $R_{i,i} = 1$ , namely its columns have unit Euclidean norm. Under this normalization, one can check that  $R$  satisfies  $(\mathcal{A}_{\text{Irr.}})$  of order  $K = 1$ . Hence, up to some normalization, one can always assume  $(\mathcal{A}_{\text{Irr.}})$  of order  $K = 1$ .

### 3.2. Conditional Joint Law of the Knots

In this section, we are interested in the joint law of the knots  $(\lambda_1, \dots, \lambda_K)$  of the LAR conditional on  $\lambda_{K+1}$ . Let  $(\widehat{i}_1, \dots, \widehat{i}_K)$  be the first variables entering along the LAR path. Define the first correlation by  $\rho_1 := \sqrt{R_{\widehat{i}_1, \widehat{i}_1}}$  and the others by

$$\rho_\ell := \frac{\sqrt{R_{\widehat{i}_\ell, \widehat{i}_\ell} - (R_{\widehat{i}_\ell, \widehat{i}_1} \cdots R_{\widehat{i}_\ell, \widehat{i}_{\ell-1}}) M_{\widehat{i}_1, \dots, \widehat{i}_{\ell-1}}^{-1} (R_{\widehat{i}_1, \widehat{i}_\ell}, \dots, R_{\widehat{i}_{\ell-1}, \widehat{i}_\ell})}}{1 - \theta_{\widehat{i}_\ell}^{\ell-1}}, \quad \text{for } \ell \geq 2, \quad (13)$$

where

$$\theta^{\ell-1} := \theta(\widehat{i}_1, \dots, \widehat{i}_{\ell-1}), \quad \text{for } \ell \geq 2,$$

is defined by (11) and  $M_{\widehat{i}_1, \dots, \widehat{i}_{\ell-1}}$  is the submatrix of  $R$  obtained by keeping the columns and the rows indexed by  $\{\widehat{i}_1, \dots, \widehat{i}_{\ell-1}\}$ . Furthermore, denote

$$F_i := \Phi_i(\lambda_i) \quad \text{and} \quad \mathcal{P}_{i,j} := \Phi_i \circ \Phi_j^{-1}, \quad \text{for } i, j \in [K+1], \quad (14)$$

where  $\Phi_k(\cdot) := \Phi(\cdot / (\sigma \rho_k))$  is the CDF of the centered Gaussian law with variance  $\sigma^2 \rho_k^2$  for  $k \geq 1$ ,  $\lambda_0 = \infty$  and  $F_0 = 1$  by convention. The main discovery of this paper is the following theorem.

**Theorem 3** (Conditional Joint Law of the LAR Knots). *Let  $(\lambda_1, \dots, \lambda_K, \lambda_{K+1})$  be the first knots and let  $(\widehat{i}_1, \dots, \widehat{i}_K)$  be the first variables entering along the LAR path. If  $(\widehat{\mathcal{A}}_{\text{irr}})$  holds then, conditional on  $\{\widehat{i}_1, \dots, \widehat{i}_K, \lambda_{K+1}\}$ , the vector  $(\lambda_1, \dots, \lambda_K)$  has law with density (w.r.t. the Lebesgue measure)*

$$Z_{(\widehat{i}_1, \dots, \widehat{i}_K, \lambda_{K+1})}^{-1} \left( \prod_{k=1}^K \varphi_{m_k, v_k^2}(\ell_k) \right) \mathbf{1}_{\{\ell_1 \geq \ell_2 \geq \dots \geq \ell_K \geq \lambda_{K+1}\}} \text{ at point } (\ell_1, \ell_2, \dots, \ell_K),$$

where  $\varphi_{m_k, v_k^2}$  is the standard Gaussian density with mean

$$m_k := \frac{\mu_{\widehat{i}_k}^0 - (R_{\widehat{i}_k, \widehat{i}_1} \cdots R_{\widehat{i}_k, \widehat{i}_{k-1}}) M_{\widehat{i}_1, \dots, \widehat{i}_{k-1}}^{-1} (\mu_{\widehat{i}_1}^0, \dots, \mu_{\widehat{i}_{k-1}}^0)}{1 - \theta_{\widehat{i}_k}^{k-1}} \quad \text{with} \quad \mu^0 := (\mu^0, -\mu^0), \quad (15)$$

and variance  $v_k^2 := \sigma^2 \rho_k^2$ .

A useful corollary of this theorem is the following. It shows that one can explicitly describe the joint law of the LAR's knots after having selected a support of size  $\widehat{m}$  with any procedure satisfying  $(\mathcal{P}_1)$  and  $(\mathcal{P}_2)$ .

**Corollary 4.** *Under the conditions of Theorem 3, let  $\widehat{m}$  be chosen according to a procedure satisfying  $(\mathcal{P}_1)$  and  $(\mathcal{P}_2)$  with  $n_{\text{select}} \geq K$ . Then under the null hypothesis, namely*

$$\mathbf{H}_0 : "X\beta^0 \in E_{\widehat{m}}",$$

*and conditional on the selection event  $\{\widehat{m} = a, F_a, F_{K+1}, \widehat{i}_1, \dots, \widehat{i}_K, \widehat{i}_{K+1}\}$  with  $0 \leq a \leq K-1$ , the vector  $(F_{a+1}, \dots, F_K)$  is uniformly distributed on*

$$\mathcal{D}_{a+1, K} := \{(f_{a+1}, \dots, f_K) \in \mathbb{R}^{K-a} : \mathcal{P}_{a+1, a}(F_a) \geq f_{a+1} \geq \mathcal{P}_{a+1, a+2}(f_{a+2}) \geq \dots \geq \mathcal{P}_{a+1, K}(f_K) \geq \mathcal{P}_{a+1, K+1}(F_{K+1})\},$$

where  $\mathcal{P}_{i,j}$  are described in (14).

**Remark 8.** *The previous statement is consistent with the case  $a = 0$  corresponding to the global null hypothesis  $\mathbf{H}_0 : "X\beta^0 = 0"$  (or equivalently  $\mathbb{E}Z = 0$ ). Therefore, if  $Z$  is centered then, conditional on  $F_{K+1}$ ,  $(F_1, \dots, F_K)$  is uniformly distributed on*

$$\mathcal{D}_{1, K} := \{(f_1, \dots, f_K) \in \mathbb{R}^K : 1 \geq f_1 \geq \mathcal{P}_{1, 2}(f_2) \geq \dots \geq \mathcal{P}_{1, K}(f_K) \geq \mathcal{P}_{1, K+1}(F_{K+1})\}.$$

**Remark 9.** In the orthogonal case where  $\bar{R} = \text{Id}$ , note that  $\theta_j(i_1, \dots, i_\ell) = 0$  for all  $\ell \geq 1$  and all  $i_1, \dots, i_\ell \neq j$ ,  $\rho_j = 1$  and  $\mathcal{P}_{i,j}(f) = f$ . We recover that  $\mathcal{D}_{1,K}$  is the set of order statistics

$$1 \geq f_1 \geq f_2 \geq \dots \geq f_K \geq \Phi(\lambda_{K+1}).$$

In this case, knots  $\lambda_i$  are Gaussian order statistics  $\lambda_1 = Z_{\hat{i}_1} \geq \lambda_2 = Z_{\hat{i}_2} \geq \dots \geq \lambda_K = Z_{\hat{i}_K} \geq \lambda_{K+1}$  of the vector  $Z$ .

### 3.3. Testing Procedures

From Theorem 3, we deduce several testing statistics. To this end, we introduce some notation. First, define

$$\mathcal{I}_{ab}(s, t) := \int_{\mathcal{P}_{(a+1),b}(t)}^{\mathcal{P}_{(a+1),a}(s)} df_{a+1} \int_{\mathcal{P}_{(a+2),b}(t)}^{\mathcal{P}_{(a+2),(a+1)}(f_{a+1})} df_{a+2} \int_{\mathcal{P}_{(a+3),b}(t)}^{\mathcal{P}_{(a+3),(a+2)}(f_{a+2})} df_{a+3} \cdots \int_{\mathcal{P}_{(b-1),b}(t)}^{\mathcal{P}_{(b-1),(b-2)}(f_{b-2})} df_{b-1}$$

for  $0 \leq a < b$  and  $s, t \in \mathbb{R}$ , with the convention that  $\mathcal{I}_{ab} = 1$  when  $b = a + 1$ ,

and also

$$\mathbb{F}_{abc}(t) := \mathbb{1}_{\{\lambda_c \leq t \leq \lambda_a\}} \int_{\Phi_b(\lambda_c)}^{\Phi_b(t)} \mathcal{I}_{ab}(F_a, f_b) \mathcal{I}_{bc}(f_b, F_c) df_b \quad (16)$$

for  $0 \leq a < b < c \leq K + 1$ ,  $t \in \mathbb{R}$  where  $F_a = \Phi_a(\lambda_a)$  and  $F_c = \Phi_c(\lambda_c)$ .

Then, a simple integration shows that Corollary 4 implies that

$$\mathbb{P}[\lambda_b \leq t \mid \lambda_a, \lambda_c, \hat{i}_1, \hat{i}_2, \dots, \hat{i}_{c-1}] = \frac{\mathbb{F}_{abc}(t)}{\mathbb{F}_{abc}(\lambda_a)}, \quad (17)$$

whenever  $\mathbb{E}Z = 0$ . Actually, this result can be refined in the case when  $\mathbb{E}Z \neq 0$  by the next result.

**Theorem 5.** Let  $(\lambda_1, \dots, \lambda_K, \lambda_{K+1})$  be the first knots and let  $(\hat{i}_1, \dots, \hat{i}_K)$  be the first variables entering along the LAR path. If  $(\hat{\mathcal{A}}_{\text{irr}})$  holds and  $\hat{m}$  is chosen according to a procedure satisfying  $(\mathcal{P}_1)$  and  $(\mathcal{P}_2)$  with  $n_{\text{select}} \geq K$ , then under the null hypothesis, namely

$$\mathbb{H}_0 : "X\beta^0 \in E_{\hat{m}}",$$

and conditional on the selection event  $\{\hat{m} = a\}$  with  $a \leq K - 1$ , and for any integers  $b, c$  such that  $a < b < c \leq K + 1$ , it holds that

$$\hat{\alpha}_{abc} := 1 - \frac{\mathbb{F}_{abc}(\lambda_b)}{\mathbb{F}_{abc}(\lambda_a)} \sim \mathcal{U}(0, 1), \quad (18)$$

namely, it is uniformly distributed over  $(0, 1)$ .

This testing statistic generalizes previous testing statistics that appeared in “Spacing Tests”, as presented in [29, Chapter 5] for instance, and will be referred to as the *Generalized Spacing test*.

**Remark 10.** If one considers  $a = 0$ ,  $b = 1$  and  $c = 2$  then one gets

$$\hat{\alpha}_{012} = 1 - \frac{\Phi_1(\lambda_1) - \Phi_1(\lambda_2)}{\Phi_1(\lambda_0) - \Phi_1(\lambda_2)} = \frac{1 - \Phi_1(\lambda_1)}{1 - \Phi_1(\lambda_2)}.$$

Similarly, taking  $b = a + 1$  and  $c = a + 2$  one gets

$$\hat{\alpha}_{a(a+1)(a+2)} = \frac{\Phi_{a+1}(\lambda_{a+1}) - \Phi_{a+1}(\lambda_a)}{\Phi_{a+1}(\lambda_{a+2}) - \Phi_{a+1}(\lambda_a)}.$$

which is the spacing test as presented in [29, Chapter 5].

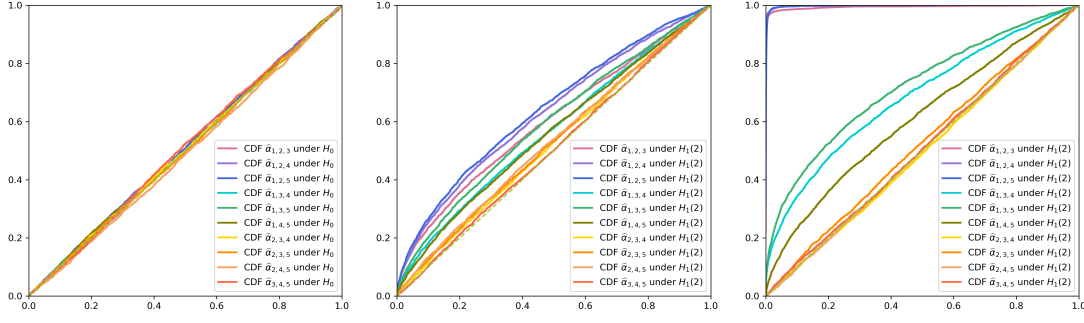


Figure 3. CDF of  $p$ -values  $\hat{\alpha}_{abc}$  over 3,000 Monte-Carlo iterations with  $n = 200$ ,  $p = 300$  and a random design matrix  $X$  given by 300 independent column vectors uniformly distributed on the Euclidean sphere  $\mathbb{S}^{299}$ . Central panel represents an alternative composed by a 2 sparse vector, the right panel an alternative composed by a 2 sparse vector 5 times larger while the left panel corresponds to the null.

Given  $\alpha \in (0, 1)$ , one can consider the following testing procedures

$$\mathcal{S}_{abc} := \mathbb{1}_{\{\hat{\alpha}_{abc} \leq \alpha\}}, \quad (19)$$

that rejects if the  $p$ -value  $\hat{\alpha}_{abc}$  is less than the level  $\alpha$  of the test. Now, recall (17) to witness that the smaller  $\hat{\alpha}_{abc}$ , the bigger  $\lambda_b$  conditionally to  $(\lambda_a, \lambda_c)$ . So **the  $p$ -value  $\hat{\alpha}_{abc}$  detects abnormal large values of  $\lambda_b$  conditional on  $(\lambda_a, \lambda_c)$ .**

One may investigate the power of these tests to detect false negatives, namely to detect the alternative given by: there exists  $k \in S^0$  such that  $k \notin \{\bar{i}_1, \dots, \bar{i}_a\}$ . In particular, *what is the most powerful test among these latter testing procedures?* A comprehensive study of the orthogonal case is given in the following theorem.

**Theorem 6.** Assume that the design is orthogonal, namely  $\bar{R} = \text{Id}_p$ . Let  $\hat{m}$  be chosen according to a procedure satisfying  $(\mathcal{P}_1)$  and  $(\mathcal{P}_2)$  with  $n_{\text{select}} \geq K$ . Then under the null hypothesis, namely

$$\mathbb{H}_0 : "X\beta^0 \in E_{\hat{m}}",$$

and conditional on the selection event  $\{\hat{m} = a_0\}$  with  $1 \leq a_0 \leq K-1$ , it holds the test  $\mathcal{S}_{a_0, a_0+1, K+1}$  is uniformly most powerful than any of the tests  $\mathcal{S}_{a,b,c}$  for  $a_0 \leq a < b < c \leq K+1$ .

The proof of this result is given in Appendix C.6. It shows that the best choice among the tests  $\mathcal{S}_{a,b,c}$  is the test  $\mathcal{S}_{a_0, a_0+1, K+1}$  with the smallest  $a$  and the largest  $c$ .

### 3.4. “Studentization”: Adapting to Unknown Noise Level

The previous results can be further refined to the case where  $\sigma^2$  is unknown. Let  $K$  chosen following the considerations of Section 2.5, the key property is to build an estimation of the variance  $\hat{\sigma}_{\text{test}}^2$  on the strata  $(n_{\text{test}}, n_{\text{select}}]$  which is independent from

$$(\lambda_{\hat{m}}, \dots, \lambda_{K+1}, \hat{\sigma}_{\text{test}}),$$

conditional on the selection event  $\{\hat{i}_1 = i_1, \dots, \hat{i}_n = i_n\}$  and under  $\mathbb{H}_0 : "X\beta^0 \in E_{\hat{m}}"$ .

Because of our hypotheses,  $\text{rank}(P_{(n_{\text{test}}, n_{\text{select}}]})$  is maximal equal to  $m := n_{\text{select}} - n_{\text{test}}$ . Our variance estimator

$$\hat{\sigma}^2 := \hat{\sigma}_{\text{test}}^2 = \frac{\|P_{(n_{\text{test}}, n_{\text{select}}]}(Z)\|_2^2}{m},$$

follows a  $\sigma^2 \chi^2(m)/m =: \sigma^2 W^2$ .

We set  $\widehat{\lambda}_i := \frac{\lambda_i}{\sigma} = \frac{\lambda_i}{\sigma W}$ . Let us consider  $a = m < b < c \leq K + 1$  and consider, under  $\mathbb{H}_0$ , the distribution conditional to  $\{\widehat{m} = a, \lambda_a, \lambda_{K+1}, \widehat{i}_1, \dots, \widehat{i}_K, \widehat{i}_{K+1}\}$ . Then the variables  $(\lambda_a, \dots, \lambda_c)$  and  $W$  are independent. More over the distribution of  $\lambda_b$  is given by (17). We have

$$\mathbb{P}\{\widehat{\lambda}_b \leq t | \widehat{\lambda}_a \widehat{\lambda}_c, W = w, \widehat{i}_1, \dots, \widehat{i}_{K+1}\} = \mathbb{P}\left\{\frac{\lambda_b}{\sigma} \leq wt | \widehat{\lambda}_a \widehat{\lambda}_c, W = w, \widehat{i}_1, \dots, \widehat{i}_{K+1}\right\}.$$

Then we can use Theorem 5 in the case  $\sigma = 1$  to get that the expression above is equal to

$$\frac{\overline{F}_{abc}(tw)}{\overline{F}_{abc}(\widehat{\lambda}_a w)},$$

Where  $\overline{F}_{abc}$  is the function given by (17) in the case  $\sigma = 1$ . De-conditioning in  $w$  we get

$$\mathbb{P}\{\widehat{\lambda}_b \leq t | \widehat{\lambda}_a \widehat{\lambda}_c, i_1, \dots, i_{K+1}\} = \int_{\mathbb{R}} \frac{\overline{F}_{abc}(tw)}{\overline{F}_{abc}(\widehat{\lambda}_a w)} d\mathbb{P}_W(w).$$

As a consequence the  $p$ -value of the test is now given by

$$\widehat{\beta}_{abc} = 1 - \int_{\mathbb{R}} \frac{\overline{F}_{abc}(\widehat{\lambda}_b w)}{\overline{F}_{abc}(\widehat{\lambda}_a w)} d\mathbb{P}_W(w).$$

### 3.5. Examples of Admissible procedures

We are now able to present here an admissible procedure to build an estimate  $\widehat{S}$  of the support that satisfies the properties  $(\mathcal{P}_1)$  and  $(\mathcal{P}_2)$ . We chose a level  $\alpha'$  and we define a light modification of the Student test

$$T_{abc} = \mathbf{1}_{\widehat{\beta}_{abc} \leq \alpha'},$$

by replacing  $\widehat{\sigma}_{\text{test}}^2$  by  $\widehat{\sigma}_{\text{select}}^2$ . The number of degrees of freedom of the  $\chi^2$  distribution is now  $\overline{m} := n - n_{\text{select}}$ . By a small abuse of notation, we still set

$$\widehat{\lambda}_i := \frac{\lambda_i}{\widehat{\sigma}_{\text{select}}}.$$

We limit our attention to consecutive  $a, b, c$ . In such a case it is easy to see that

$$\widehat{\beta}_{a(a+1)(a+2)} = \frac{\mathcal{T}\left(\frac{\widehat{\lambda}_{a+1}}{\sigma \rho_{a+1}}\right) - \mathcal{T}\left(\frac{\widehat{\lambda}_a}{\sigma \rho_{a+1}}\right)}{\mathcal{T}\left(\frac{\widehat{\lambda}_{a+2}}{\sigma \rho_{a+1}}\right) - \mathcal{T}\left(\frac{\widehat{\lambda}_a}{\sigma \rho_{a+1}}\right)},$$

where  $\mathcal{T}$  is the cumulative distribution of the Student  $T(\overline{m})$  law, so the quantity above is easy to compute. We are now in condition to present our algorithm.

- Begin with  $a = 0$ ,
- at each step, perform the test  $T_{a,a+1,a+2}$  at the level  $\alpha'$ ,
- if the test is significative, we set  $a = a + 1$  and keep on going,
- if it is non-significative, we stop and set  $\widehat{m} = a + 2$ .

Recall that possible selected supports along the LAR's path are nested models of the form (2). Denote  $k^0 \geq 1$  the smallest integer  $k$  such that the true support  $S^0$  is contained in  $\overline{S}^k$  and denote

$$S^0 \subseteq \overline{S}^{k^0} \text{ and } S^0 \not\subseteq \overline{S}^{(k^0-1)}.$$

We understand that admissible procedures depend on the sequence of false positives appearing along the LAR's path. If the experimenter believes that there is no more than  $\gamma_{FP}$  **consecutive** false positives in  $S^0$  a possible admissible procedure would be the following.



- Begin with  $a = 0$ ,
- at each step, perform the test  $T_{a,a+1,a+2}$  at the level  $\alpha'$ ,
- if the test is significative, we set  $a = a + 1$  and keep on going,
- if the  $\gamma_{FP}$  consecutive tests  $T_{a,a+1,a+2}, T_{a+1,a+2,a+3}, \dots, T_{a+\gamma_{FP}-1,a+\gamma_{FP},a+\gamma_{FP}+1}$  are all non-significatives, we stop and set  $\hat{m} = a + \gamma_{FP} + 1$ .

This method has been deployed on real data in Section 4.2 with  $\gamma_{FP} = 3$ .

## 4. Numerical Experiments

### 4.1. Monte-Carlo experiment

To study the relative power in a *non-orthogonal* design case we have build a Monte-Carlo experiment with 3,000 repetitions. We have considered a model with  $n = 200$ ,  $p = 300$  and a random design matrix  $X$  given by 300 independent column vectors uniformly distributed on the Euclidean sphere  $\mathbb{S}^{299}$ . The computation of the function  $F_{abc}$  given by (16) is an important issue that demands multivariate integration tools, see Appendix E for a solution using cubature of integral by lattice rule. This has lead to some limitations namely  $\Delta_{\text{num}} \leq 4$  that implies  $c \leq 5$  when  $a = 1$  in our experimental framework.

A python notebook and codes are given at [https://github.com/ydecastro/lar\\_testing](https://github.com/ydecastro/lar_testing). The base function is `observed_significance_CBC(lars, sigma, start, end, middle)` in the file `multiple_spacing_tests.py`. It gives the  $p$ -value of  $T_{(\text{start})(\text{middle})(\text{end})}$  of knots and indexes given by `lars` and an estimation of (or the true) standard deviation given by `sigma`. We have run 3,000 repetitions of this function to get the laws displayed in Figure 3. It presents the CDF of the  $p$ -value  $\hat{\alpha}_{abc}$  under the null and under two 2-sparse alternatives, one with low signal and one with 5 times more signal. Results show, in our particular case, that all the tests are exact and the test  $\mathcal{S}_{125}$  is the most powerful. More precisely, it holds, as in the orthogonal case, that

- $\mathcal{S}_{123} \leq \mathcal{S}_{124} \leq \mathcal{S}_{125}$ ;
- $\mathcal{S}_{234} \leq \mathcal{S}_{245}$ .

### 4.2. Real data

A detailed presentation in a Python notebook is available at [https://github.com/ydecastro/lar\\_testing/blob/master/multiple\\_spacing\\_tests.ipynb](https://github.com/ydecastro/lar_testing/blob/master/multiple_spacing_tests.ipynb).

We consider a data set about HIV drug resistance extracted from [3] and [22]. The experiment consists in identifying mutations on the genes of the HIV-virus that are involved with drug resistance. The data set contains about  $p = 200$  and  $n = 700$  observations. Since some protocol is used to remove some gene or some individuals, the exact numbers depend on the considered drug.

We used a procedure referred to as “*spacing-BH*” procedure which is a Benjamini–Hochberg procedure based on the sequence of spacing tests

$$\hat{\beta}_{012}, \hat{\beta}_{123}, \dots, \hat{\beta}_{a(a+1)(a+2)}, \dots$$

as described in Section 3.4 with  $\alpha = 0.2$ . The results for Knockoff of [3] and of Benjamini–Hochberg procedure on the coefficients of linear regression (BHq) are for the R-vignette `knockoff` of the dedicated web page of Stanford. All results are evaluated using the TSM data base that gives, in some sense, the list of true positives. A comparison of our results with those of [3] is displayed in Figure 4. Our procedure is a bit more conservative but, in most of the case, gives a better control of the FDP.

In addition we have performed on the same dataset a false negative detection as in Section 2.4, and Section 3.5 with  $\gamma_{FP} = 3$ . We refer to the aforementioned Python notebook for further details.

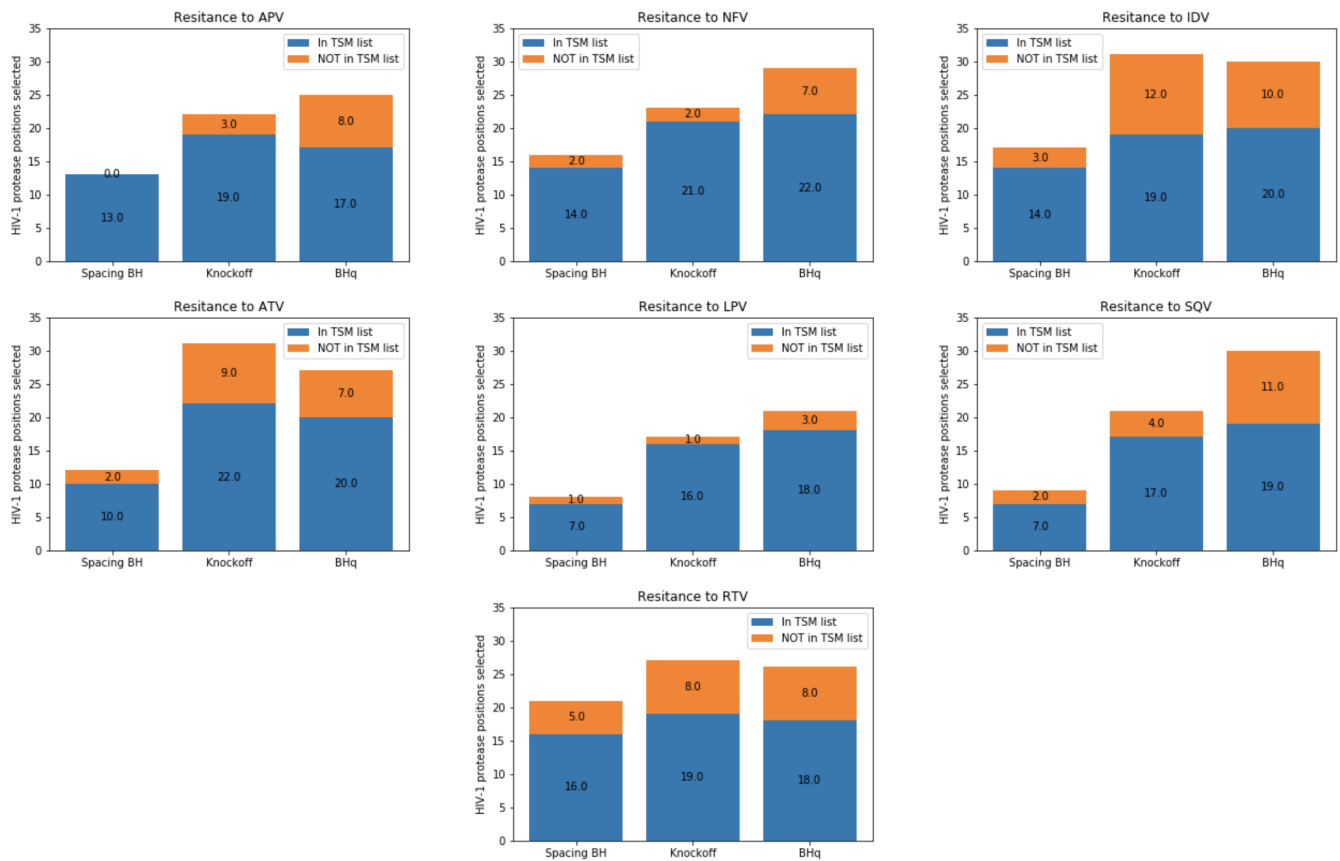


Figure 4. Comparison of numbers true and false positives for three procedures: Spacing-BH; Knockoff [3] and Benjamini-Hochberg procedure on the coefficient of linear regression BHq. For each drug we indicate the number of true positives in blue and false positives in orange. In the three procedures the aimed FDR is  $\alpha = 20\%$ .

## References

- [1] J.-M. Azais, Y. De Castro, and S. Mourareau. Power of the spacing test for least-angle regression. *Bernoulli*, 24(1):465–492, 2018.
- [2] F. Bachoc, G. Blanchard, P. Neuvial, et al. On the post selection inference constant under restricted isometry properties. *Electronic Journal of Statistics*, 12(2):3736–3757, 2018.
- [3] R. F. Barber, E. J. Candès, et al. Controlling the false discovery rate via knockoffs. *The Annals of Statistics*, 43(5):2055–2085, 2015.
- [4] P. C. Bellec, G. Lecué, A. B. Tsybakov, et al. Slope meets lasso: improved oracle bounds and optimality. *The Annals of Statistics*, 46(6B):3603–3642, 2018.
- [5] Y. Benjamini and Y. Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*, 57(1):289–300, 1995.
- [6] R. Berk, L. Brown, A. Buja, K. Zhang, L. Zhao, et al. Valid post-selection inference. *The Annals of Statistics*, 41(2):802–837, 2013.
- [7] P. J. Bickel, Y. Ritov, A. B. Tsybakov, et al. Simultaneous analysis of lasso and dantzig selector. *The Annals of Statistics*, 37(4):1705–1732, 2009.
- [8] G. Blanchard, P. Neuvial, and E. Roquain. Post hoc inference via joint family-wise error rate control. *arXiv preprint arXiv:1703.02307*, 2017.
- [9] G. Blanchard, E. Roquain, et al. Two simple sufficient conditions for fdr control. *Electronic journal of Statistics*, 2:963–992, 2008.
- [10] M. Bogdan, E. Van Den Berg, C. Sabatti, W. Su, and E. J. Candès. Slope—adaptive variable selection via convex optimization. *The annals of applied statistics*, 9(3):1103, 2015.
- [11] P. Bühlmann and S. van de Geer. *Statistics for high-dimensional data*. Springer Series in Statistics. Springer, Heidelberg, 2011. Methods, theory and applications.
- [12] E. J. Candès, J. Romberg, and T. Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Trans. Inf. Theory*, 52(2):489–509, 2006.
- [13] S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM J. Sci. Comput.*, 20(1):33–61 (electronic), 1998.
- [14] B. Efron, T. Hastie, I. Johnstone, R. Tibshirani, et al. Least angle regression. *The Annals of statistics*, 32(2):407–499, 2004.
- [15] W. Fithian, D. Sun, and J. Taylor. Optimal inference after model selection. *arXiv preprint arXiv:1410.2597*, 2014.
- [16] A. Genz. Numerical computation of multivariate normal probabilities. *Journal of computational and graphical statistics*, 1(2):141–149, 1992.
- [17] C. Giraud. *Introduction to high-dimensional statistics*. Chapman and Hall/CRC, 2014.
- [18] A. Javanmard, H. Javadi, et al. False discovery rate control via debiased lasso. *Electronic Journal of Statistics*, 13(1):1212–1253, 2019.
- [19] F. Le Gall. Powers of tensors and fast matrix multiplication. In *Proceedings of the 39th International Symposium on Symbolic and Algebraic Computation*, ISSAC ’14, pages 296–303, New York, NY, USA, 2014. ACM.
- [20] R. Lockhart, J. Taylor, R. J. Tibshirani, and R. Tibshirani. A significance test for the lasso. *Annals of statistics*, 42(2):413, 2014.
- [21] D. Nuyens and R. Cools. Fast algorithms for component-by-component construction of rank-1 lattice rules in shift-invariant reproducing kernel hilbert spaces. *Mathematics of Computation*, 75(254):903–920, 2006.
- [22] S.-Y. Rhee, J. Taylor, G. Wadhera, A. Ben-Hur, D. L. Brutlag, and R. W. Shafer. Genotypic predictors of human immunodeficiency virus type 1 drug resistance. *Proceedings of the National Academy of Sciences*, 103(46):17355–17360, 2006.
- [23] E. Roquain. Type i error rate control for testing many hypotheses: a survey with proofs. *Journal de la Société Française de Statistique*, 152(2):3–38, 2011.
- [24] F. Santambrogio. Optimal transport for applied mathematicians. *Birkhäuser*, NY, 55:58–63, 2015.

- [25] J. Taylor and R. J. Tibshirani. Statistical learning and selective inference. *Proceedings of the National Academy of Sciences*, 112(25):7629–7634, 2015.
- [26] X. Tian, J. R. Loftus, and J. E. Taylor. Selective inference with unknown variance via the square-root lasso. *Biometrika*, 105(4):755–768, 2018.
- [27] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- [28] R. Tibshirani, M. Wainwright, and T. Hastie. *Statistical Learning with Sparsity: The Lasso and Generalizations*. Monographs on Statistics & Applied Probability. Chapman and Hall/CRC press, 2015.
- [29] R. J. Tibshirani, J. Taylor, R. Lockhart, and R. Tibshirani. Exact post-selection inference for sequential regression procedures. *Journal of the American Statistical Association*, 111(514):600–620, 2016.
- [30] S. van de Geer. Estimation and testing under sparsity. *Lecture Notes in Mathematics*, 2159, 2016.
- [31] C. Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008.
- [32] M. J. Wainwright. Sharp thresholds for high-dimensional and noisy sparsity recovery using  $\ell_1$ -constrained quadratic programming (lasso). *IEEE transactions on information theory*, 55(5):2183–2202, 2009.

## Appendix A: Representing the LAR Knots

### A.1. The equivalent formulations of the LAR algorithm

We present here three equivalent formulations of the LAR that are a consequence of the analysis provided in Appendices A and B. One formulation is given by Algorithm 1.

---

#### Algorithm 2: LAR algorithm (standard formulation)

---

**Data:** Correlations vector  $\bar{Z}$  and variance-covariance matrix  $\bar{R}$ .

**Result:** Sequence  $((\lambda_k, \bar{v}_k, \varepsilon_k))_{k \geq 1}$  where  $\lambda_1 \geq \lambda_2 \geq \dots > 0$  are the knots, and  $\bar{v}_1, \bar{v}_2, \dots$  are the variables that enter the model with signs  $\varepsilon_1, \varepsilon_2, \dots$  ( $\varepsilon_k = \pm 1$ ).

*/\* Initialize computing  $(\lambda_1, \bar{v}_1, \varepsilon_1)$  and defining a “residual”  $\bar{N}^{(1)}$ . \*/*

1 Set  $k = 1$ ,  $\lambda_1 := \max |\bar{Z}|$ ,  $\bar{v}_1 := \arg \max |\bar{Z}|$  and  $\varepsilon_1 = \bar{Z}_{\bar{v}_1} / \lambda_1 \in \pm 1$ , and  $\bar{N}^{(1)} := \bar{Z}$ .

*/\* Note that  $((\lambda_\ell, \bar{v}_\ell, \varepsilon_\ell))_{1 \leq \ell \leq k-1}$  and  $\bar{N}^{(k-1)}$  have been defined at the previous step. \*/*

2 Set  $k \leftarrow k + 1$  and compute the least-squares fit

$$\bar{\theta}_j := (\bar{R}_{j, \bar{v}_1} \cdots \bar{R}_{j, \bar{v}_{k-1}}) M_{\bar{v}_1, \dots, \bar{v}_{k-1}}^{-1} (\varepsilon_1, \dots, \varepsilon_{k-1}), \quad j = 1, \dots, p,$$

where  $M_{\bar{v}_1, \dots, \bar{v}_{k-1}}$  is the submatrix of  $\bar{R}$  keeping the columns and the rows indexed by  $\{\bar{v}_1, \dots, \bar{v}_{k-1}\}$ .

3 For  $0 < \lambda \leq \lambda_{k-1}$  compute the “residuals”  $\bar{N}^{(k)}(\lambda) = (\bar{N}_1^{(k)}(\lambda), \dots, \bar{N}_p^{(k)}(\lambda))$  given by

$$\bar{N}_j^{(k)}(\lambda) := \bar{N}_j^{(k-1)} - (\lambda_{k-1} - \lambda) \bar{\theta}_j, \quad j = 1, \dots, p,$$

and pick

$$\lambda_k := \max \{ \beta > 0; \exists j \notin \{\bar{v}_1, \dots, \bar{v}_{k-1}\}, \text{ s.t. } |\bar{N}_j^{(k)}(\beta)| = \beta \} \text{ and } \bar{v}_k := \arg \max_{j \notin \{\bar{v}_1, \dots, \bar{v}_{k-1}\}} |\bar{N}_j^{(k)}(\lambda_k)|,$$

$$\varepsilon_k := \bar{N}_{\bar{v}_k}^{(k)}(\lambda_k) / \lambda_k \in \pm 1 \text{ and } \bar{N}^{(k)} := \bar{N}^{(k)}(\lambda_k).$$

Then, iterate from 2.

---



---

#### Algorithm 3: LAR algorithm (“projected” formulation)

---

**Data:** Correlations vector  $\bar{Z}$  and variance-covariance matrix  $\bar{R}$ .

**Result:** Sequence  $((\lambda_k, \bar{v}_k, \varepsilon_k))_{k \geq 1}$  where  $\lambda_1 \geq \lambda_2 \geq \dots > 0$  are the knots, and  $\bar{v}_1, \bar{v}_2, \dots$  are the variables that enter the model with signs  $\varepsilon_1, \varepsilon_2, \dots$  ( $\varepsilon_k = \pm 1$ ).

*/\* Initialize computing  $(\lambda_1, \bar{v}_1, \varepsilon_1)$ . \*/*

1 Define  $Z = (\bar{Z}, -\bar{Z})$  and  $R$  as in (10), and set  $k = 1$ ,  $\lambda_1 := \max Z$ ,  $\hat{v}_1 := \arg \max Z$ ,  $\bar{v}_1 = \hat{v}_1 \bmod p$  and  $\varepsilon_1 = 1 - 2(\hat{v}_1 - \bar{v}_1)/p \in \pm 1$ .

*/\* Note that  $((\lambda_\ell, \hat{v}_\ell))_{1 \leq \ell \leq k-1}$  have been defined at the previous step/loop. \*/*

2 Set  $k \leftarrow k + 1$  and compute

$$\lambda_k = \max_{\{j: \theta_j(\hat{v}_1, \dots, \hat{v}_{k-1}) < 1\}} \left\{ \frac{Z_j - P_{\hat{v}_1, \dots, \hat{v}_{k-1}}(Z_j)}{1 - \theta_j(\hat{v}_1, \dots, \hat{v}_{k-1})} \right\} \text{ and } \hat{v}_k = \arg \max_{\{j: \theta_j(\hat{v}_1, \dots, \hat{v}_{k-1}) < 1\}} \left\{ \frac{Z_j - P_{\hat{v}_1, \dots, \hat{v}_{k-1}}(Z_j)}{1 - \theta_j(\hat{v}_1, \dots, \hat{v}_{k-1})} \right\},$$

where

$$P_{\hat{v}_1, \dots, \hat{v}_{k-1}}(Z_j) := (R_{j, \hat{v}_1} \cdots R_{j, \hat{v}_{k-1}}) M_{\hat{v}_1, \dots, \hat{v}_{k-1}}^{-1} (Z_{\hat{v}_1}, \dots, Z_{\hat{v}_{k-1}})$$

$$\theta_j(\hat{v}_1, \dots, \hat{v}_{k-1}) := (R_{j, \hat{v}_1} \cdots R_{j, \hat{v}_{k-1}}) M_{\hat{v}_1, \dots, \hat{v}_{k-1}}^{-1} (1, \dots, 1)$$

and set  $\bar{v}_k = \hat{v}_k \bmod p$  and  $\varepsilon_k = 1 - 2(\hat{v}_k - \bar{v}_k)/p \in \pm 1$ . Then, iterate from 2.

---

### A.2. Initialization: First Knot

The first step of the LAR algorithm (Step 1 in Algorithm 2) seeks the most correlated predictor with the observation. In our formulation, introduce the first residual  $N^{(1)} := Z$  and observe that  $N^{(1)} := (\bar{N}^{(1)}, -\bar{N}^{(1)})$ . We define the first knot  $\lambda_1 > 0$  as

$$\lambda_1 = \max Z \quad \text{and} \quad \hat{v}_1 = \arg \max Z.$$

One may see that this definition is consistent with  $\lambda_1$  in Algorithm 2 and note that  $\hat{v}_1$  and  $(\bar{v}_1, \varepsilon_1)$  are related as in (1).

The LAR algorithm is a forward algorithm that selects a new variable and maintains a residual at each step. We also define

$$N^{(2)}(\lambda) = N^{(1)} - (\lambda_1 - \lambda)\theta(\hat{v}_1), \quad 0 < \lambda \leq \lambda_1, \quad (20)$$

and one can check that  $N^{(2)}(\lambda) = (\bar{N}^{(2)}(\lambda), -\bar{N}^{(2)}(\lambda))$  where  $\bar{N}(\lambda)$  is defined in Algorithm 2. It is clear that the coordinate  $\hat{v}_1$  of  $N^{(2)}(\lambda)$  is equal to  $\lambda$ . On the other hand  $N^{(1)} = Z$  attains its maximum at the single point  $\hat{v}_1$ . By continuity this last property is kept for  $\lambda$  in a left neighborhood of  $\lambda_1$ . We search for the first value of  $\lambda$  such that this property is not met, *i.e.* the largest value of  $\lambda$  such that

$$\exists j \neq \hat{v}_1 \text{ such that } N^{(2)}(\lambda) = \lambda,$$

as in Step 3 of Algorithm 2. We call this value  $\lambda_2$  and one may check that this definition is consistent with  $\lambda_2$  in Algorithm 2.

Now, we can be more explicit on the expression of  $\lambda_2$ . Indeed, we may make the following discussion depending on the values of  $\theta_j(\hat{v}_1)$ .

- If  $\theta_j(\hat{v}_1) \geq 1$ , since  $N_j^{(1)} < N_{\hat{v}_1}^{(1)}$  for  $j \neq \hat{v}_1$  there is no hope to achieve the equality between  $N_j^{(2)}(\lambda)$  and  $N_{\hat{v}_1}^{(2)}(\lambda) = \lambda$  for  $0 < \lambda \leq \lambda_1$  in view of (20).
- Thus we limit our attention to the  $j$ 's such that  $\theta_j(\hat{v}_1) < 1$ . We have equality  $N_j^{(2)}(\lambda) = \lambda$  when

$$\lambda = \frac{N_j^{(1)} - \lambda_1 \theta_j(\hat{v}_1)}{1 - \theta_j(\hat{v}_1)}.$$

So we can also define the second knot  $\lambda_2$  of the LAR as

$$\lambda_2 = \max_{j: \theta_j(\hat{v}_1) < 1} \left\{ \frac{Z_j - P_{\hat{v}_1}(Z_j)}{1 - \theta_j(\hat{v}_1)} \right\}.$$

where  $P_{i_1}(Z_j) := Z_{i_1} \theta_j(i_1)$ . Remark that  $P_{i_1}(Z_j) = \mathbb{E}(Z_j \mid Z_{i_1})$  is the regression of  $Z_j$  on  $Z_{i_1}$  when  $\mathbb{E}Z = 0$ .

### A.3. Recursion: Next Knots

The loop (2  $\Leftrightarrow$  3) in Algorithm 2 builds iteratively the knots  $\lambda_1, \lambda_2, \dots$  of the LAR algorithm and some “residuals”  $\bar{N}^{(1)}, \bar{N}^{(2)}, \dots$  defined in Step 3. We will present here an equivalent formulation of these knots.

Assume that  $k \geq 2$  and we have build  $\lambda_1, \dots, \lambda_{k-1}$  and selected the “signed” variables  $\hat{v}_1, \dots, \hat{v}_{k-1}$ . Introduce  $N^{(k-1)} := (\bar{N}^{(k-1)}, -\bar{N}^{(k-1)})$  and define

$$N^{(k)}(\lambda) = N^{(k-1)} - (\lambda_{k-1} - \lambda)\theta(\hat{v}_1, \dots, \hat{v}_{k-1}), \quad 0 < \lambda \leq \lambda_{k-1}.$$

Check that  $\theta_j(\hat{v}_1, \dots, \hat{v}_{k-1}) = (\bar{\theta}_j, -\bar{\theta}_j)$  where we recall that we define

$$\bar{\theta}_j := (\bar{R}_{j, \bar{v}_1} \cdots \bar{R}_{j, \bar{v}_{k-1}}) M_{\bar{v}_1, \dots, \bar{v}_{k-1}}^{-1} (\varepsilon_1, \dots, \varepsilon_{k-1}), \quad j = 1, \dots, p,$$

at Step 2 and it holds that  $\widehat{v}_\ell$  and  $(\widehat{v}_\ell, \varepsilon_\ell)$  are related as in (1). From this equality, we deduce that it holds  $N^{(k)}(\lambda) = (\overline{N}^{(k)}(\lambda), -\overline{N}^{(k)}(\lambda))$ . One may also check that the coordinates  $\widehat{v}_1, \dots, \widehat{v}_{k-1}$  of  $N^{(k)}(\lambda)$  are equal to  $\lambda$ .

Again if we want to solve  $N_j^{(k)}(\lambda) = \lambda$  for some  $j$ , we have to limit our attention to  $j$ 's such that  $\theta_j(\widehat{v}_1, \dots, \widehat{v}_{k-1}) < 1$ . Solving this latter equality yields to

$$\lambda_k = \max_{j: \theta_j(\widehat{v}_1, \dots, \widehat{v}_{k-1}) < 1} \left\{ \frac{N_j^{(k-1)} - \lambda_{k-1} \theta_j(\widehat{v}_1, \dots, \widehat{v}_{k-1})}{1 - \theta_j(\widehat{v}_1, \dots, \widehat{v}_{k-1})} \right\}.$$

This expression is consistent with  $\lambda_k$  in Algorithm 2.

Now, we can give an other expression of  $\lambda_k$  that will be useful in the proofs of our main theorems. Note that the residuals satisfy the relation

$$N^{(k)} = N^{(k-1)} - (\lambda_{k-1} - \lambda_k) \theta(\widehat{v}_1, \dots, \widehat{v}_{k-1}), \quad (21)$$

and that  $N_j^{(k-1)} = \lambda_{k-1}$  for  $j = \widehat{v}_1, \dots, \widehat{v}_{k-1}$ . The following lemma permits a drastic simplification of the expression of the knots. Its proof is given in Appendix C.3.

**Lemma 7.** *It holds*

$$N^{(k-1)} - \lambda_{k-1} \theta(\widehat{v}_1, \dots, \widehat{v}_{k-1}) = Z - P_{\widehat{v}_1, \dots, \widehat{v}_{k-1}}(Z)$$

where we denote  $P_{i_1, \dots, i_{k-1}}(Z) = (P_{i_1, \dots, i_{k-1}}(Z_1), \dots, P_{i_1, \dots, i_{k-1}}(Z_{2p}))$  and for all  $j \in [2p]$  one has  $P_{i_1, \dots, i_{k-1}}(Z_j) = (R_{j, i_1} \cdots R_{j, i_{k-1}}) M_{i_1, \dots, i_{k-1}}^{-1}(Z_{i_1}, \dots, Z_{i_{k-1}})$ .

Using Lemma 7 we deduce that  $\lambda_k$  in Algorithm 2 is consistent with

$$\lambda_k = \max_{j: \theta_j(\widehat{v}_1, \dots, \widehat{v}_{k-1}) < 1} \left\{ \frac{Z_j - P_{\widehat{v}_1, \dots, \widehat{v}_{k-1}}(Z_j)}{1 - \theta_j(\widehat{v}_1, \dots, \widehat{v}_{k-1})} \right\}.$$

where  $P_{\widehat{v}_1, \dots, \widehat{v}_{k-1}}(Z_j) = (R_{j, \widehat{v}_1} \cdots R_{j, \widehat{v}_{k-1}}) M_{\widehat{v}_1, \dots, \widehat{v}_{k-1}}^{-1}(Z_{\widehat{v}_1}, \dots, Z_{\widehat{v}_{k-1}})$ . When  $\mathbb{E}Z = 0$ , one may remark that  $P_{i_1, \dots, i_{k-1}}(Z_j)$  is the regression of  $Z_j$  on the vector  $(Z_{i_1}, \dots, Z_{i_{k-1}})$  whose variance-covariance matrix is  $M_{i_1, \dots, i_{k-1}}$ . This analysis leads to an equivalent formulation of the LAR algorithm (Algorithm 2). We present this formulation in Algorithm 3.

**Remark 11.** *Note that Algorithm 2 implies that  $\widehat{v}_1, \dots, \widehat{v}_k$  are pairwise different, but also that they differ modulo  $p$ . In the rest of the paper we will limit our attention to such sequences.*

## Appendix B: First Steps to Derive the Joint Law of the LAR Knots

Given fixed  $i_1, \dots, i_k \in [2p]$ , one may define  $Z_j^{(i_1, \dots, i_k)} := Z_j - P_{i_1, \dots, i_k}(Z_j)$  for indices  $j$  such that  $\theta_j(i_1, \dots, i_k) = 1$  and

$$\forall j \text{ s.t. } \theta_j(i_1, \dots, i_k) \neq 1, \quad Z_j^{(i_1, \dots, i_k)} := \frac{Z_j - P_{i_1, \dots, i_k}(Z_j)}{1 - \theta_j(i_1, \dots, i_k)},$$

where we recall that  $P_{i_1, \dots, i_k}(Z_j) = (R_{j, i_1} \cdots R_{j, i_k}) M_{i_1, \dots, i_k}^{-1}(Z_{i_1}, \dots, Z_{i_k})$ . From this point, we can introduce

$$\lambda_{k+1}^{(i_1, \dots, i_k)} := \max_{j: \theta_j(i_1, \dots, i_k) < 1} Z_j^{(i_1, \dots, i_k)},$$

and remark that  $\lambda_{k+1} = \lambda_{k+1}^{(\widehat{v}_1, \dots, \widehat{v}_k)}$ .



### B.1. Law of the First Knot

One has the following lemma governing the law of  $\lambda_1$ .

**Lemma 8.** *It holds that*

- $Z_{i_1}$  is independent of  $(Z_j^{(i_1)})_{j \neq i_1}$ ,
- If  $\theta_j(i_1) < 1$  for all  $j \neq i_1$  then

$$\{\widehat{i}_1 = i_1\} = \{\lambda_2^{(i_1)} \leq Z_{i_1}\},$$

- If  $\theta_j(i_1) < 1$  for all  $j \neq i_1$  then, conditionally to  $\{\widehat{i}_1 = i_1\}$  and  $\lambda_2, \lambda_1$  is a truncated Gaussian random variable with mean  $\mathbb{E}(Z_{i_1})$  and variance  $\rho_1^2 := R_{\widehat{i}_1, \widehat{i}_1}$  subject to be greater than  $\lambda_2$ .

*Proof.* The first point is clear and standard. Now, observe that

$$\begin{aligned} \{\lambda_2^{(i_1)} \leq Z_{i_1}\} &\Leftrightarrow \{\forall j \neq i_1, \frac{Z_j - Z_{i_1}\theta_j(i_1)}{1 - \theta_j(i_1)} \leq Z_{i_1}\} \\ &\Leftrightarrow \{\forall j \neq i_1, Z_j - Z_{i_1}\theta_j(i_1) \leq Z_{i_1} - Z_{i_1}\theta_j(i_1)\} \\ &\Leftrightarrow \{\forall j \neq i_1, Z_j \leq Z_{i_1}\} \\ &\Leftrightarrow \{\widehat{i}_1 = i_1\}, \end{aligned}$$

as claimed. The last statement is a consequence of the two previous points.  $\square$

### B.2. Recursive Formulation of the LAR

One has the following proposition whose proof can be found in Section C.4. As we will see in this section, this intermediate result as a deep consequence, the LAR algorithm can be stated in a recursive way applying the same function repeatedly, as presented in Algorithm 1.

**Proposition 9.** *Set*

$$\tau_{j, i_k} := \frac{R_{j, i_k} - (R_{j, i_1} \cdots R_{j, i_{k-1}}) M_{i_1, \dots, i_{k-1}}^{-1} (R_{i_k, i_1}, \dots, R_{i_k, i_{k-1}})}{(1 - \theta_j(i_1, \dots, i_{k-1}))(1 - \theta_{i_k}(i_1, \dots, i_{k-1}))},$$

and observe that  $\tau_{j, i_k}$  is the covariance between  $Z_j^{(i_1, \dots, i_{k-1})}$  and  $Z_{i_k}^{(i_1, \dots, i_{k-1})}$ . Furthermore, it holds

$$\frac{\tau_{j, i_k}}{\tau_{i_k, i_k}} = 1 - \frac{1 - \theta_j(i_1, \dots, i_k)}{1 - \theta_{i_k}(i_1, \dots, i_{k-1})} \quad (22)$$

and

$$\forall j \neq i_1, \dots, i_k, \quad Z_j^{(i_1, \dots, i_k)} = \frac{Z_j^{(i_1, \dots, i_{k-1})} - Z_{i_k}^{(i_1, \dots, i_{k-1})} \tau_{j, i_k} / \tau_{i_k, i_k}}{1 - \tau_{j, i_k} / \tau_{i_k, i_k}}. \quad (23)$$

Now, we present Algorithm 1. Define  $R(0) := R$ ,  $Z(0) = Z$  and  $T(0) = 0$ . For  $k \geq 1$  and fixed  $i_1, \dots, i_k \in [2p]$ , introduce

$$\begin{aligned} R(k) &:= \left( R_{j, \ell} - (R_{j, i_1} \cdots R_{j, i_k}) M_{i_1, \dots, i_k}^{-1} (R_{\ell, i_1}, \dots, R_{\ell, i_k}) \right)_{j, \ell} \\ Z(k) &:= Z - P_{i_1, \dots, i_k}(Z) \\ T(k) &:= (\theta_j(i_1, \dots, i_k))_j, \end{aligned}$$

and note that  $R(k)$  is the variance-covariance matrix of the Gaussian vector  $Z(k)$ . The key property is following. Let  $v_1, \dots, v_k$ , be  $k$  linearly independent vectors of an Euclidean space and let  $u$  be any vector of the space. Set

$$v := P_{(v_1, \dots, v_{k-1})}^\perp v_k,$$

the projection of  $v_k$  orthogonally to  $v_1, \dots, v_{k-1}$ . Then

$$P_{(v_1, \dots, v_k)}^\perp u = P_v^\perp P_{(v_1, \dots, v_{k-1})} u.$$

Using this result we deduce that

$$\begin{aligned} Z(k) &= P_{i_1, \dots, i_k}^\perp(Z) \\ &= P_{i_k}^\perp(P_{i_1, \dots, i_{k-1}}^\perp(Z)) \\ &= P_{i_k}^\perp(Z(k-1)) \\ &= Z(k-1) - P_{i_k}(Z(k-1)) \\ &= Z(k-1) - \mathbf{x}(k)Z(k-1), \end{aligned} \tag{24}$$

where  $\mathbf{x}(k) = R_{i_k}(k-1)/R_{i_k, i_k}(k-1)$ . It yields that

$$R(k) = R(k-1) - \mathbf{x}(k)R_{i_k}(k-1)^\top. \tag{25}$$

Using (22) (or (31)), remark that

$$T(k) = T(k-1) - \mathbf{x}(k)(1 - T_{i_k}(k-1)). \tag{26}$$

These relations give a recursive formulation of the LAR as presented in Algorithm 1.

### B.3. Joint Law

Regarding the joint law of knots, one has the following proposition whose proof can be found in Section C.5.

**Proposition 10.** *One has the following for any fixed  $i_1, \dots, i_{k+1} \in [2p]$ .*

- *It holds that*

$$(Z_j^{(i_1, \dots, i_{k+1})})_{j \neq i_1, \dots, i_{k+1}} \perp\!\!\!\perp Z_{i_{k+1}}^{(i_1, \dots, i_k)} \perp\!\!\!\perp Z_{i_k}^{(i_1, \dots, i_{k-1})} \perp\!\!\!\perp \dots \perp\!\!\!\perp Z_{i_2}^{(i_1)} \perp\!\!\!\perp Z_{i_1}$$

*are mutually independent.*

- *If  $(\hat{\mathcal{A}}_{\text{Irr.}})$  of order  $k$  holds then*

$$\begin{aligned} &\{\hat{v}_1 = i_1, \dots, \hat{v}_{k+1} = i_{k+1}\} \\ &= \{\lambda_{k+1}^{(i_1, \dots, i_k)} = Z_{i_{k+1}}^{(i_1, \dots, i_k)} \leq Z_{i_k}^{(i_1, \dots, i_{k-1})} \leq \dots \leq Z_{i_2}^{(i_1)} \leq Z_{i_1}\}, \\ &= \{\lambda_{k+1}^{(i_1, \dots, i_k)} = Z_{i_{k+1}}^{(i_1, \dots, i_k)} \leq Z_{i_k}^{(i_1, \dots, i_{k-1})} \leq \dots \leq Z_{i_{a+1}}^{(i_1, \dots, i_a)} \leq Z_{i_a}^{(i_1, \dots, i_{a-1})} = \lambda_a^{(i_1, \dots, i_{a-1})}\} \\ &\quad \bigcap \{\hat{v}_1 = i_1, \dots, \hat{v}_a = i_a, \hat{v}_{k+1} = i_{k+1}\} \bigcap \{\lambda_{K+1} = Z_{i_{k+1}}^{(i_1, \dots, i_k)}, \dots, \lambda_a = Z_{i_a}^{(i_1, \dots, i_{a-1})}\}, \end{aligned}$$

*for any  $0 \leq a \leq k-1$  with the convention  $\lambda_0 = \infty$ .*

## Appendix C: Proofs

### C.1. Proof of Theorem 3

#### Step 1

We prove by induction on  $k$  the following equality

$$\{\widehat{i}_1 = i_1, \dots, \widehat{i}_{K+1} = i_{K+1}\} = \left\{ \lambda_{k+1}^{(i_1, \dots, i_k)} \leq \left[ \begin{array}{c} \lambda_k^{(i_1, \dots, i_{k-1})} \\ Z_{i_k}^{(i_1, \dots, i_{k-1})} \end{array} \leq \dots \leq \left[ \begin{array}{c} \lambda_1 \\ Z_{i_1} \end{array} \right] \right\},$$

where the bracket means that, at any position you can choose indifferently the upper- or the lower-case to get an equivalent event.

Using Lemma 8 we have the following equivalences

$$\{\widehat{i}_1 = i_1\} = \{\lambda_1 = Z_{i_1}\} = \{\lambda_2^{(i_1)} \leq Z_{i_1}\} = \{\lambda_2^{(i_1)} \leq \lambda_1\}$$

that give the property for  $k = 1$ .

Using the induction hypothesis and Proposition 9 we can pass from  $k$  to  $k + 1$  to get that

$$\{\widehat{i}_1 = i_1, \dots, \widehat{i}_{k+1} = i_{k+1}\} = \left\{ \lambda_{k+2}^{(i_1, \dots, i_{k+1})} \leq \lambda_{k+1}^{(i_1, \dots, i_k)} \leq \left[ \begin{array}{c} \lambda_k^{(i_1, \dots, i_{k-1})} \\ Z_{i_k}^{(i_1, \dots, i_{k-1})} \end{array} \leq \dots \leq \left[ \begin{array}{c} \lambda_1 \\ Z_{i_1} \end{array} \right] \right\}.$$

The  $2^k$  equivalent events above imply obviously the following  $2^k$  events

$$\left\{ \lambda_{k+2}^{(i_1, \dots, i_{k+1})} \leq Z_{i_{k+1}}^{(i_1, \dots, i_k)} \leq \left[ \begin{array}{c} \lambda_k^{(i_1, \dots, i_{k-1})} \\ Z_{i_k}^{(i_1, \dots, i_{k-1})} \end{array} \leq \dots \leq \left[ \begin{array}{c} \lambda_1 \\ Z_{i_1} \end{array} \right] \right\}.$$

To get the implication in the other direction we remark that

$$\begin{aligned} \lambda_{k+2}^{(i_1, \dots, i_{k+1})} &\leq Z_{i_{k+1}}^{(i_1, \dots, i_k)} \\ \Leftrightarrow \forall j \neq i_1, \dots, i_{k+1}, \quad Z_j^{(i_1, \dots, i_{k+1})} &\leq Z_{i_{k+1}}^{(i_1, \dots, i_k)} \\ \Leftrightarrow \forall j \neq i_1, \dots, i_{k+1}, \quad Z_j^{(i_1, \dots, i_k)} - Z_{i_{k+1}}^{(i_1, \dots, i_k)} \frac{\tau_{j, i_{k+1}}}{\tau_{i_{k+1}, i_{k+1}}} &\leq Z_{i_{k+1}}^{(i_1, \dots, i_k)} - Z_{i_{k+1}}^{(i_1, \dots, i_k)} \frac{\tau_{j, i_{k+1}}}{\tau_{i_{k+1}, i_{k+1}}} \\ \Leftrightarrow \forall j \neq i_1, \dots, i_{k+1}, \quad Z_j^{(i_1, \dots, i_k)} &\leq Z_{i_{k+1}}^{(i_1, \dots, i_k)} \\ \Leftrightarrow \lambda_{k+1}^{(i_1, \dots, i_k)} &= Z_{i_{k+1}}^{(i_1, \dots, i_k)}. \end{aligned} \tag{27}$$

using that (23) and that  $1 - \tau_{j, i_{k+1}} / \tau_{i_{k+1}, i_{k+1}} > 0$  (which is a consequence of (22) and  $(\widehat{\mathcal{A}}_{\text{Irr.}})$  in (27).

#### Step 2

Using the version 3 of the LAR algorithm we deduce directly the expressions of the expectation and variance of

$$Z_{i_\ell}^{(i_1, \dots, i_{\ell-1})}$$

given by (13) (15) (Changing  $k$  into  $\ell$ ). From this same version of this algorithm we get the independence of the  $Z_{i_\ell}^{(i_1, \dots, i_{\ell-1})}$  when  $\ell$  varies.

Now suppose that we are conditional to  $(a, \lambda_a, \lambda_{k+1}, \widehat{i}_1, \dots, \widehat{i}_{K+1})$ , then standard properties of order statistics imply that the density of  $(\lambda_{a+1}, \dots, \lambda_k)$  taken at  $(\ell_{a+1}, \dots, \ell_k)$  is proportional to

$$\left( \prod_{k=a+1}^K \varphi_{m_k, v_k^2}(\ell_k) \right) \mathbb{1}_{\{\ell_{a+1} \geq \ell_{a+2} \geq \dots \geq \ell_K \geq \lambda_{K+1}\}}. \tag{28}$$

### Step 3

Suppose now that we are conditional to  $\{\hat{m} = a, \lambda_a, \lambda_{K+1}, \hat{v}_1, \dots, \hat{v}_K, \hat{v}_{K+1}\}$  with  $0 \leq a \leq K-1$ . Because of  $\mathcal{P}_2$ ,  $\lambda_{a+1}, \dots, \lambda_K$  are independent of the condition  $m = a$ , so their conditional distribution is equal to the unconditional distribution given by (28).

### C.2. Proof of Proposition 2

Let  $S = \{j_1, \dots, j_k\} \subset [p]$  and  $j \in [2p] \setminus S$ . Let  $\bar{v} = (\bar{v}_1, \dots, \bar{v}_k) \in \{-1, 1\}^k$  and define  $i_\ell = j_\ell + p(1 - \bar{v}_\ell)/2$  for  $\ell \in [k]$ . Note that

$$\begin{aligned} \theta_j(i_1, \dots, i_k) &= (R_{j,i_1} \cdots R_{j,i_k}) M_{i_1, \dots, i_k}^{-1} (1, \dots, 1) = \left[ X_j^\top X_S \text{Diag}(\bar{v}) \right] M_{i_1, \dots, i_k}^{-1} (1, \dots, 1) \\ &= \left[ X_j^\top X_S \text{Diag}(\bar{v}) \right] M_{i_1, \dots, i_k}^{-1} \left[ \text{Diag}(\bar{v}) \bar{v} \right] = X_j^\top X_S \left[ \text{Diag}(\bar{v}) M_{i_1, \dots, i_k}^{-1} \text{Diag}(\bar{v}) \right] \bar{v} \\ &= X_j^\top X_S (X_S^\top X_S)^{-1} \bar{v}. \end{aligned}$$

Now, observe that

$$\max_{\|v\|_\infty \leq 1} v^\top (X_S^\top X_S)^{-1} X_S^\top X_j = \max_{\bar{v} \in \{-1, 1\}^k} X_j^\top X_S (X_S^\top X_S)^{-1} \bar{v},$$

showing the equivalence between the two assumptions.

### C.3. Proof of Lemma 7

The proof works by induction. Let us check the relation for  $k = 2$ , namely

$$N^{(1)} - \lambda_1 \theta(\hat{v}_1) = Z - Z_{\hat{v}_1} \theta(\hat{v}_1) = Z - P_{\hat{v}_1}(Z).$$

Now, let  $k \geq 3$ . First, the three perpendicular theorem implies that for every  $j, i_1, \dots, i_{k-1}$ ,

$$\begin{aligned} \theta_j(i_1, \dots, i_{k-2}) &= (R_{j,i_1} \cdots R_{j,i_{k-1}}) M_{i_1, \dots, i_{k-1}}^{-1} (\theta_{i_1}(i_1, \dots, i_{k-2}), \dots, \theta_{i_{k-1}}(i_1, \dots, i_{k-2})), \\ \text{and } P_{i_1, \dots, i_{k-2}}(Z_j) &= (R_{j,i_1} \cdots R_{j,i_{k-1}}) M_{i_1, \dots, i_{k-1}}^{-1} (P_{i_1, \dots, i_{k-2}}(Z_{i_1}), \dots, P_{i_1, \dots, i_{k-2}}(Z_{i_{k-1}})). \end{aligned}$$

By induction, using (21), we get that

$$\begin{aligned} N^{(k-1)} &= N^{(k-2)} - (\lambda_{k-2} - \lambda_{k-1}) \theta(\hat{v}_1, \dots, \hat{v}_{k-2}), \\ &= (N^{(k-2)} - \lambda_{k-2} \theta(\hat{v}_1, \dots, \hat{v}_{k-2})) + \lambda_{k-1} \theta(\hat{v}_1, \dots, \hat{v}_{k-2}), \\ &= Z - P_{\hat{v}_1, \dots, \hat{v}_{k-2}}(Z) + \lambda_{k-1} \theta(\hat{v}_1, \dots, \hat{v}_{k-2}). \end{aligned} \tag{29}$$

Then, recall that  $N_j^{(k-1)} = \lambda_{k-1}$  for  $j = \hat{v}_1, \dots, \hat{v}_{k-1}$  and remark that

$$\lambda_{k-1} \theta_j(\hat{v}_1, \dots, \hat{v}_{k-1}) = (R_{j,\hat{v}_1} \cdots R_{j,\hat{v}_{k-1}}) M_{\hat{v}_1, \dots, \hat{v}_{k-1}}^{-1} (N_{\hat{v}_1}^{(k-1)}, \dots, N_{\hat{v}_{k-1}}^{(k-1)}).$$

Using (29) at indices  $j = \hat{v}_1, \dots, \hat{v}_{k-1}$ , we deduce that

$$\begin{aligned} \lambda_{k-1} \theta_j(\hat{v}_1, \dots, \hat{v}_{k-1}) &= (R_{j,\hat{v}_1} \cdots R_{j,\hat{v}_{k-1}}) M_{\hat{v}_1, \dots, \hat{v}_{k-1}}^{-1} (N_{\hat{v}_1}^{(k-1)}, \dots, N_{\hat{v}_{k-1}}^{(k-1)}) \\ &= (R_{j,\hat{v}_1} \cdots R_{j,\hat{v}_{k-1}}) M_{\hat{v}_1, \dots, \hat{v}_{k-1}}^{-1} (Z_{\hat{v}_1}, \dots, Z_{\hat{v}_{k-1}}) \\ &\quad - (R_{j,\hat{v}_1} \cdots R_{j,\hat{v}_{k-1}}) M_{\hat{v}_1, \dots, \hat{v}_{k-1}}^{-1} (P_{\hat{v}_1, \dots, \hat{v}_{k-2}}(Z_{\hat{v}_1}), \dots, P_{\hat{v}_1, \dots, \hat{v}_{k-2}}(Z_{\hat{v}_{k-1}})) \\ &\quad + \lambda_{k-1} (R_{j,\hat{v}_1} \cdots R_{j,\hat{v}_{k-1}}) M_{\hat{v}_1, \dots, \hat{v}_{k-1}}^{-1} (\theta_{\hat{v}_1}(\hat{v}_1, \dots, \hat{v}_{k-2}), \dots, \theta_{\hat{v}_{k-1}}(\hat{v}_1, \dots, \hat{v}_{k-2})) \\ &= P_{\hat{v}_1, \dots, \hat{v}_{k-1}}(Z_j) - P_{\hat{v}_1, \dots, \hat{v}_{k-2}}(Z_j) + \lambda_{k-1} \theta_j(\hat{v}_1, \dots, \hat{v}_{k-2}), \end{aligned}$$

Namely

$$P_{\widehat{i}_1, \dots, \widehat{i}_{k-1}}(Z) - \lambda_{k-1} \theta(\widehat{i}_1, \dots, \widehat{i}_{k-1}) = P_{\widehat{i}_1, \dots, \widehat{i}_{k-2}}(Z) - \lambda_{k-1} \theta(\widehat{i}_1, \dots, \widehat{i}_{k-2}).$$

Using again (29) we get that

$$\begin{aligned} N^{(k-1)} &= Z - P_{\widehat{i}_1, \dots, \widehat{i}_{k-2}}(Z) + \lambda_{k-1} \theta(\widehat{i}_1, \dots, \widehat{i}_{k-2}), \\ &= Z - P_{\widehat{i}_1, \dots, \widehat{i}_{k-1}}(Z) + \lambda_{k-1} \theta(\widehat{i}_1, \dots, \widehat{i}_{k-1}), \end{aligned}$$

as claimed.

#### C.4. Proof of Proposition 9

We denote

$$\begin{aligned} R_j &:= (R_{j, i_1}, \dots, R_{j, i_{k-1}}), \\ R_{i_k} &:= (R_{i_k, i_1}, \dots, R_{i_k, i_{k-1}}), \\ M &:= M_{i_1, \dots, i_{k-1}}, \\ \overline{M} &:= M_{i_1, \dots, i_k} = \begin{bmatrix} M & R_{i_k} \\ R_{i_k}^\top & R_{i_k, i_k} \end{bmatrix}, \\ \overline{R} &:= (R_{j, i_1}, \dots, R_{j, i_k}), \\ x &:= \frac{1 - \theta_j(i_1, \dots, i_{k-1})}{1 - \theta_{i_k}(i_1, \dots, i_{k-1})} \frac{\tau_{j, i_k}}{\tau_{i_k, i_k}}, \end{aligned}$$

and observe that

$$\begin{aligned} x &= \frac{R_{j, i_k} - R_j^\top M^{-1} R_{i_k}}{R_{i_k, i_k} - R_{i_k}^\top M^{-1} R_{i_k}}, \\ \overline{M}^{-1} &= \begin{bmatrix} \text{Id}_{k-1} & -M^{-1} R_{i_k} \\ 0 & 1 \end{bmatrix} \begin{bmatrix} M^{-1} & 0 \\ 0 & (R_{i_k, i_k} - R_{i_k}^\top M^{-1} R_{i_k})^{-1} \end{bmatrix} \begin{bmatrix} \text{Id}_{k-1} & 0 \\ -R_{i_k}^\top M^{-1} & 1 \end{bmatrix}, \\ \overline{M}^{-1} \overline{R} &= \begin{bmatrix} M^{-1} (R_j - x R_{i_k}) \\ x \end{bmatrix}, \end{aligned} \tag{30}$$

using the Schur complement of the block  $M$  of the matrix  $\overline{M}$  and a LU decomposition. Note also that

$$\frac{Z_j^{(i_1, \dots, i_{k-1})} - Z_{i_k}^{(i_1, \dots, i_{k-1})} \tau_{j, i_k} / \tau_{i_k, i_k}}{1 - \tau_{j, i_k} / \tau_{i_k, i_k}} = \frac{Z_j - P_{i_1, \dots, i_{k-1}}(Z_j) - x (Z_{i_k} - P_{i_1, \dots, i_{k-1}}(Z_{i_k}))}{1 - \theta_j(i_1, \dots, i_{k-1}) - x (1 - \theta_{i_k}(i_1, \dots, i_{k-1}))}.$$

To prove (23), it suffices to show that the R.H.S term above is equal to the following R.H.S term

$$Z_j^{(i_1, \dots, i_k)} = \frac{Z_j - P_{i_1, \dots, i_k}(Z_j)}{1 - \theta_j(i_1, \dots, i_k)}.$$

We will prove that the numerators are equal and that the denominators are equal. For the denominators, we use that

$$\begin{aligned} &1 - \theta_j(i_1, \dots, i_{k-1}) - x (1 - \theta_{i_k}(i_1, \dots, i_{k-1})) \\ &= 1 - \theta_j(i_1, \dots, i_{k-1}) - x + x \theta_{i_k}(i_1, \dots, i_{k-1}) \\ &= 1 - (\underbrace{1 \dots 1}_{k \text{ times}}) \begin{bmatrix} M^{-1} (R_j - x R_{i_k}) \\ x \end{bmatrix} \\ &= 1 - \theta_j(i_1, \dots, i_k), \end{aligned} \tag{31}$$

using (30). Furthermore, it proves (22). For the numerators, we use that

$$\begin{aligned}
& Z_j - P_{i_1, \dots, i_{k-1}}(Z_j) - x(Z_{i_k} - P_{i_1, \dots, i_{k-1}}(Z_{i_k})) \\
&= Z_j - P_{i_1, \dots, i_{k-1}}(Z_j) - xZ_{i_k} + xP_{i_1, \dots, i_{k-1}}(Z_{i_k}) \\
&= Z_j - (Z_{i_1} \cdots Z_{i_k}) \begin{bmatrix} M^{-1}(R_j - xR_{i_k}) \\ x \end{bmatrix} \\
&= Z_j - P_{i_1, \dots, i_k}(Z_j).
\end{aligned}$$

using (30).

### C.5. Proof of Proposition 10

The proof of the first point can be lead by induction. The initialization of the proof is given by the first point of Lemma 8. Now, observe that  $Z_{i_k}^{(i_1, \dots, i_{k-1})}, \dots, Z_{i_2}^{(i_1)}, Z_{i_1}$  are measurable functions of  $(Z_{i_1}, \dots, Z_{i_k})$  and one may check that the vector  $(Z_{i_1}, \dots, Z_{i_k})$  is independent of the vector  $(Z_{i_{k+1}}^{(i_1, \dots, i_k)}, (Z_j^{(i_1, \dots, i_{k+1})})_{j \neq i_1, \dots, i_{k+1}})$ . We deduce that  $Z_{i_k}^{(i_1, \dots, i_{k-1})}, \dots, Z_{i_2}^{(i_1)}, Z_{i_1}$  are independent of  $(Z_{i_{k+1}}^{(i_1, \dots, i_k)}, (Z_j^{(i_1, \dots, i_{k+1})})_{j \neq i_1, \dots, i_{k+1}})$ . One can also check that  $Z_{i_{k+1}}^{(i_1, \dots, i_k)}$  is independent of  $(Z_j^{(i_1, \dots, i_{k+1})})_{j \neq i_1, \dots, i_{k+1}}$ .

The second point also works by induction. The initialization of the proof is given by the second point of Lemma 8. We will use Proposition 9 to prove the second point.

Now, we have

$$\begin{aligned}
\lambda_{k+1}^{(i_1, \dots, i_k)} &\leq Z_{i_k}^{(i_1, \dots, i_{k-1})} \\
&\Leftrightarrow \forall j \neq i_1, \dots, i_k, Z_j^{(i_1, \dots, i_k)} \leq Z_{i_k}^{(i_1, \dots, i_{k-1})} \\
&\Leftrightarrow \forall j \neq i_1, \dots, i_k, Z_j^{(i_1, \dots, i_{k-1})} - Z_{i_k}^{(i_1, \dots, i_{k-1})} \frac{\tau_{j, i_k}}{\tau_{i_k, i_k}} \leq Z_{i_k}^{(i_1, \dots, i_{k-1})} - Z_{i_k}^{(i_1, \dots, i_{k-1})} \frac{\tau_{j, i_k}}{\tau_{i_k, i_k}} \\
&\Leftrightarrow \forall j \neq i_1, \dots, i_k, Z_j^{(i_1, \dots, i_{k-1})} \leq Z_{i_k}^{(i_1, \dots, i_{k-1})} \\
&\Leftrightarrow \lambda_k^{(i_1, \dots, i_{k-1})} = Z_{i_k}^{(i_1, \dots, i_{k-1})}.
\end{aligned} \tag{32}$$

using that (23) and that  $1 - \tau_{j, i_k} / \tau_{i_k, i_k} > 0$  (which is a consequence of (22) and  $(\hat{\mathcal{A}}_{\text{Irr.}})$  in (32). By induction and using (32), it holds that

$$\begin{aligned}
&\{\lambda_{k+1}^{(i_1, \dots, i_k)} \leq Z_{i_k}^{(i_1, \dots, i_{k-1})} \leq Z_{i_{k-1}}^{(i_1, \dots, i_{k-2})} \leq \dots \leq Z_{i_2}^{(i_1)} \leq Z_{i_1}\} \\
&\Leftrightarrow \{\forall j \neq i_1, \dots, i_k, Z_j^{(i_1, \dots, i_{k-1})} \leq Z_{i_k}^{(i_1, \dots, i_{k-1})} \leq Z_{i_{k-1}}^{(i_1, \dots, i_{k-2})} \leq \dots \leq Z_{i_2}^{(i_1)} \leq Z_{i_1}\} \\
&\Leftrightarrow \{\forall j \neq i_1, \dots, i_{k-1}, Z_j^{(i_1, \dots, i_{k-1})} \leq Z_{i_{k-1}}^{(i_1, \dots, i_{k-2})} \leq \dots \leq Z_{i_2}^{(i_1)} \leq Z_{i_1} \\
&\quad \text{and } \forall j \neq i_1, \dots, i_k, Z_j^{(i_1, \dots, i_{k-1})} \leq Z_{i_k}^{(i_1, \dots, i_{k-1})}\} \\
&\Leftrightarrow \{\lambda_k^{(i_1, \dots, i_{k-1})} \leq Z_{i_{k-1}}^{(i_1, \dots, i_{k-2})} \leq \dots \leq Z_{i_2}^{(i_1)} \leq Z_{i_1} \text{ and } \lambda_k^{(i_1, \dots, i_{k-1})} = Z_{i_k}^{(i_1, \dots, i_{k-1})}\} \\
&\vdots \\
&\Leftrightarrow \{\lambda_a^{(i_1, \dots, i_{a-1})} \leq Z_{i_{a-1}}^{(i_1, \dots, i_{a-2})} \leq \dots \leq Z_{i_2}^{(i_1)} \leq Z_{i_1} \\
&\quad \text{and } \lambda_k^{(i_1, \dots, i_{k-1})} = Z_{i_k}^{(i_1, \dots, i_{k-1})} \leq \dots \leq \lambda_a^{(i_1, \dots, i_{a-1})} = Z_{i_a}^{(i_1, \dots, i_{a-1})}\} \\
&\vdots \\
&\Leftrightarrow \{\hat{i}_1 = i_1, \dots, \hat{i}_k = i_k\}.
\end{aligned} \tag{s_a}$$

Now, observe that  $\hat{i}_{k+1}$  is the (unique)  $\arg \max$  of  $\lambda_{k+1}^{(i_1, \dots, i_k)}$  on the event  $\{\hat{i}_1 = i_1, \dots, \hat{i}_k = i_k\}$ . It yields that

$$\{\hat{i}_1 = i_1, \dots, \hat{i}_{k+1} = i_{k+1}\} = \{\lambda_{k+1}^{(i_1, \dots, i_k)} = Z_{i_{k+1}}^{(i_1, \dots, i_k)} \leq Z_{i_k}^{(i_1, \dots, i_{k-1})} \leq \dots \leq Z_{i_2}^{(i_1)} \leq Z_{i_1}\},$$

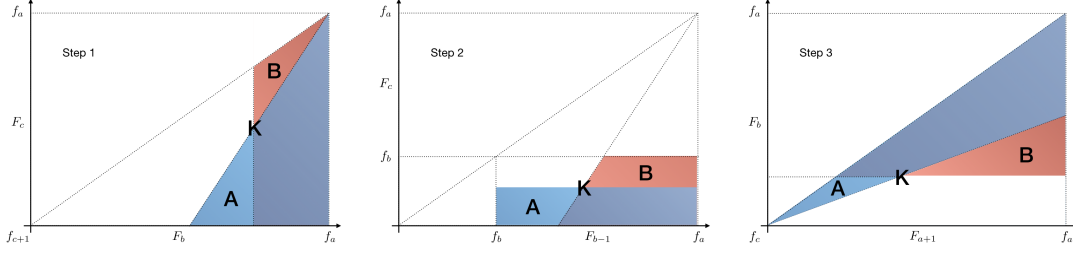


Figure 5. Rejection domains associated to the different comparison sets appearing in steps of the proof of Theorem 6.

as claimed. Stopping at  $a$  as in  $(s_a)$  gives the second part of the statement.

### C.6. Orthogonal Case: Proof of Theorem 6

Let  $\mathcal{I}$  the set of admissible indexes

$$\mathcal{I} := \{a, b, c : a_0 \leq a < b < c \leq K + 1\}.$$

◦ **Step 1:** We prove that, when the considered indexes such that  $c + 1 \leq K + 1$  belong to  $\mathcal{I}$ ,  $\mathcal{S}_{a,b,c+1}$  is more powerful than  $\mathcal{S}_{a,b,c}$ . Our proof is conditional to  $F_a = f_a, F_{c+1} = f_{c+1}$ . Note that  $(F_{a+1}, \dots, F_c)$  has for distribution the uniform distribution on the simplex

$$\mathcal{S} := \{f_a > F_{a+1} > \dots > F_c > f_{c+1}\}.$$

This implies by direct calculations that

$$\mathcal{I}_{ab}(s, t) = \frac{(s - t)^{b-a-1}}{(b - a - 1)!}$$

and that

$$\frac{\mathbb{F}_{abc}(\lambda_b)}{\mathbb{F}_{abc}(\lambda_a)} = \mathbf{F}_{\beta((b-a), (c-b))} \left( \frac{F_b - F_c}{f_a - F_c} \right). \quad (33)$$

where  $\mathbf{F}_{\beta}$  is the cumulative distribution of the Beta distribution in reference. Using monotony of this function the  $\mathcal{S}_{abc}$  test has for rejection region

$$(F_b - F_c) \geq z_1(f_a - F_c) \Leftrightarrow F_b \geq z_1 f_a + (1 - z_1) F_c, \quad (34)$$

where  $z_1$  is some threshold depending on  $\alpha$  that belongs to  $(0, 1)$ .

Similarly  $\mathcal{S}_{ab(c+1)}$  has for rejection region

$$F_b \geq z_2(f_a - f_{c+1}) + f_{c+1}, \quad (35)$$

where  $z_2$  is some other threshold belonging to  $(0, 1)$ . We use the following lemmas.

**Lemma 11.** *Let  $c \leq K$ . The density  $h_{\mu}$  of  $f_1, \dots, f_c$ , conditional to  $F_{c+1}$  with respect of the Lebesgue measure under the alternative is coordinate-wise non-decreasing and given by (36).*

*Proof.* Observe that it suffices to prove the result when  $\sigma = 1$ . Note that

$$\lambda_{c+1}^{i_1, \dots, i_c} = \max_{j \in [P], j \neq \tilde{i}_1, \dots, j \neq \tilde{i}_c} |Z_j|.$$

Thus its density  $p_{\mu^0, i_1, \dots, i_c}$  does not depend on  $\mu_{\tilde{i}_1}^0, \dots, \mu_{\tilde{i}_c}^0$ . As a consequence the following variables have the same distribution ;  $\lambda_{c+1}^{i_1 + \epsilon_1 p, \dots, i_c + \epsilon_c p}$ , where  $\epsilon_1, \dots, \epsilon_c$  take the value 0 or 1 and indices are taken modulo  $p$ .



Because of independence of the different variables, the joint density, under the alternative, of  $\lambda_1, \dots, \lambda_{c+1}$  taken at  $\ell_1, \dots, \ell_{c+1}$ , on the domain  $\{\lambda_1 > \dots > \lambda_{c+1}\}$  takes the value

$$(Const) \sum' (\phi(\ell_1 - \mu_{j_1}^0) + \phi(\ell_1 + \mu_{j_1}^0)), \dots, (\phi(\ell_c - \mu_{j_c}^0) + \phi(\ell_c + \mu_{j_c}^0)) p_{\mu^0, j_1, \dots, j_K}(\ell_{c+1}).$$

Here the sum  $\sum'$  is taken over all different  $j_1, \dots, j_c$  belonging to  $\llbracket 1, p \rrbracket$ .

Then the density, conditional to  $F_{c+1} = f_{c+1}$ , of  $F_1, \dots, F_c$  at  $f_1, \dots, f_c$  takes the value

$$(\text{const}) \sum' \cosh(\mu_{j_1} f_1) \dots \cosh(\mu_{j_c} f_c) \mathbf{1}_{f_1 > \dots > f_c > F_{c+1}}, \quad (36)$$

implying that this density is coordinate-wise non-decreasing.  $\square$

**Lemma 12.** *Let  $\nu_0$  the image on the plane  $(F_b, F_c)$  on the uniform probability on  $\mathcal{S}$ : it is the distribution under the null of  $(F_b, F_c)$ . The two rejection regions  $\mathcal{R}_1$  associated to (34) and  $\mathcal{R}_2$  associated to (35) have of course the same probability  $\alpha$  under  $\nu_0$ . Let  $\eta_{\mu^0}$  the density w.r.t.  $\nu_0$  of the distribution of  $(F_b, F_c)$  under the alternative.*

*Then  $\eta_{\mu^0}$  is non decreasing coordinate-wise.*

*Proof.* Integration yields that the density of  $\nu_0$  w.r.t. the Lebesgue measure taken at point  $(f_b, f_c)$  is

$$\frac{(f_a - f_b)^{b-a-1} (f_b - f_c)^{c-b-1}}{(b-a-1)!(c-b-1)!}.$$

The density of  $\nu_{\mu^0}$  w.r.t. Lebesgue measure is

$$\int_{f_b}^{f_a} df_{a+1} \dots \int_{f_b}^{f_b-2} df_{a+1} \int_{f_b}^{f_b-2} df_{b-1} \int_{f_c}^{f_b} df_{b+1} \dots \int_{f_c}^{f_c-2} df_{c-1} h_{\mu^0}(f_a, \dots, f_c). \quad (37)$$

Thus  $\eta_{\mu^0}$  which is the quotient of these two quantities is just a mean value of  $h_{\mu^0}$  on the domain of integration  $\mathcal{D}_{f_b, f_c}$  in (37).

Suppose that  $f_b$  and  $f_c$  increase, then all the borns of the domain  $\mathcal{D}_{f_b, f_c}$  increase also. By Lemma 11 the mean value increases.  $\square$

**We finish now the proof of Step 1:** For a given level  $\alpha$  let us consider the two rejection regions  $R_{a,b,c}$  and  $R_{a,b,(c+1)}$  of the two considered tests in the plane  $F_b, F_c$  and set

$$A := R_{a,b,c} \setminus R_{a,b,(c+1)} \text{ and } B := R_{a,b,(c+1)} \setminus R_{a,b,c},$$

see Figure 5. These two regions have the same  $\nu_0$  measure. By elementary geometry there exist a point  $K = (K_b, K_c)$  in the plane such that

- For every point of  $A$ ,  $F_b \leq K_b$ ,  $F_c \leq K_c$ ,
- For every point of  $B$ ,  $F_b \geq K_b$ ,  $F_c \geq K_c$ ,

By transport of measure there exists a transport function  $\mathcal{T}$  that preserve the measure  $\nu_0$  and that is one-to one  $A \rightarrow B$ . As a consequence the transport by  $\mathcal{T}$  improve the probability under the alternative: the power of  $\mathcal{S}_{a,b,c+1}$  is larger than that of  $\mathcal{S}_{a,b,c}$ .

◦ **Step 2:** We prove that, when the considered indexes belong to  $\mathcal{I}$  such that  $a < b-1$ ,  $\mathcal{S}_{a,(b-1),c}$  is more powerful than  $\mathcal{S}_{a,b,c}$ . Our proof is conditional to  $F_a = f_a, F_b = f_b$  and is located in the plane  $(F_{b-1}, F_c)$ .

The rejection region  $R_{a,b,c}$  takes the form  $F_c \leq \frac{1}{1-z_1} f_b - \frac{z_1}{1-z_1} f_a$  for some threshold  $z_1$  belonging to  $(0, 1)$ .

The rejection region  $R_{a,(b-1),c}$  takes the form  $F_c \leq \frac{1}{1-z_2} F_{b-1} - \frac{z_2}{1-z_2} f_a$  for some other threshold  $z_2$  belonging to  $(0, 1)$ .

These regions as well as the regions  $A$  and  $B$  and the point  $K$  are indicated in Figure 5.

Transport of measure and the convenient modification of Lemma 12 imply that the power of  $\mathcal{S}_{a,(b-1),c}$  is greater or equal than that of  $\mathcal{S}_{a,b,c}$ .

◦ **Step 3:** We prove that, when the considered indexes belong to  $\mathcal{I}$  such that  $a+1 < b$ ,  $\mathcal{S}_{a,b,c}$  is more powerful than  $\mathcal{S}_{(a+1),b,c}$ . Our proof is conditional to  $F_a = f_a, F_c = f_c$  and is located in the plane  $F_{a+1}, F_b$ .

The rejection region  $R_{a,b,c}$  takes the form  $F_b \geq z_1 f_a + (1 - z_1) f_c$  for some threshold  $z_1$  belonging to  $(0, 1)$ .

The rejection region  $R_{a+1,b,c}$  takes the form  $F_b \geq z_2 F_{a+1} + (1 - z_2) f_c$  for some other threshold  $z_2$  belonging to  $(0, 1)$ .

These regions as well as the regions  $A$  and  $B$  and the point  $K$  are indicated in Figure 5.

Transport of measure and the convenient modification of Lemma 12 imply that the power of  $\mathcal{S}_{a,b,c}$  is greater or equal than that of  $\mathcal{S}_{(a+1),b,c}$ .

Considering the three cases above, we get the desired result.

### C.7. Proof of Theorem 1

We rely on the **Weak Positive Regression Dependency (WPRDS) property** to prove the result, one may consult [17, Page 173] for instance. We say that a function  $g : [0, 1]^K \rightarrow \mathbb{R}^+$  is *nondecreasing* if for any  $p, q \in [0, 1]^K$  such that  $p_k \geq q_k$  for any  $k = 1, \dots, K$ , we have  $g(p) \geq g(q)$ . We say that a Borel set  $\Gamma \in [0, 1]^K$  is *nondecreasing* if  $g = \mathbb{1}_\Gamma$  is nondecreasing. In other words if  $y \in \gamma$  and if  $z \geq 0$ , then  $y + z \in \gamma$ . We say that the  $p$ -values  $(\hat{p}_1 = \hat{\alpha}_{0,1,K+1}, \dots, \hat{p}_K = \hat{\alpha}_{K-1,K,K+1})$  satisfy the WPRDS if for any nondecreasing set  $\Gamma$  and for all  $k^0 \in I_0$ , the function

$$u \mapsto \mathbb{P}_{\mu^0}[(\hat{p}_1, \dots, \hat{p}_K) \in \Gamma | \hat{p}_{k^0} \leq u] \text{ is nondecreasing}$$

where  $\mu^0 = \beta^0$  in our orthogonal design case, and we recall that

$$I_0 = \{k \in [K] : \mathbb{H}_{0,k} \text{ is true}\}.$$

To prove Theorem 1, note that it is sufficient [17, Chapter 8] to prove that

$$u \mapsto \overline{\mathbb{P}}[(\hat{p}_1, \dots, \hat{p}_K) \in \Gamma | \hat{p}_{k^0} \leq u] \text{ is nondecreasing} \quad (38)$$

where  $\overline{\mathbb{E}}, \overline{\mathbb{P}}$  will denote that expectations and probabilities are conditional on  $\{\bar{\iota}_1, \dots, \bar{\iota}_K, \lambda_{K+1}\}$  and under the hypothesis that  $\mu^0 = X^\top X \beta^0$ . Note that one can integrate in  $\lambda_{K+1}$  to get the statement of Theorem 1.

◦ **Step 1:** We start by giving the joint law of the LAR's knots under the alternative in the orthogonal design case. Lemma 11 and (36) show that, conditional on  $\{\bar{\iota}_1, \dots, \bar{\iota}_K, \lambda_{K+1}\}$ ,  $(\lambda_1, \dots, \lambda_K)$  is distributed on the set  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_K \geq \lambda_{K+1}$  and it has a coordinate-wise nondecreasing density. Now we can assume without loss of generality that  $\sigma^2 = 1$ , in addition because of orthogonality  $\rho_k^2 = 1$  implying that  $F_k = \Phi(\lambda_k)$   $\mathcal{P}_{i,j} = \Phi_i \circ \Phi_j^{-1} = \text{Id}$ . We deduce that, conditional on  $\{\bar{\iota}_1, \dots, \bar{\iota}_K, F_{K+1}\}$ ,  $(F_1, \dots, F_K)$  is distributed on the set

$$\{(f_1, \dots, f_K) \in \mathbb{R}^K : 1 \geq f_1 \geq f_2 \geq \dots \geq f_K \geq F_{K+1}\},$$

it has an **explicit density** given by (36), and we denote it by  $h_{\mu^0}$ . By the change of variables  $G_k := \frac{F_k - F_{K+1}}{F_{k-1} - F_{K+1}}$  one obtains that the distribution of  $(G_1, \dots, G_K)$  is supported on  $[0, 1]^K$ . More precisely, define

$$\begin{aligned} \varphi(f_1, \dots, f_K) &:= (g_1, \dots, g_K) := \left( \frac{f_1 - F_{K+1}}{1 - F_{K+1}}, \dots, \frac{f_K - F_{K+1}}{f_{K-1} - F_{K+1}} \right) \\ \varphi^{-1}(g_1, \dots, g_K) &:= \left( (1 - F_{K+1})g_1 + F_{K+1}, \dots, (1 - F_{K+1})g_K + F_{K+1} \right), \end{aligned}$$

whose inverse Jacobian determinant is

$$\det \left[ \frac{\partial \psi}{\partial f_1} \cdots \frac{\partial \psi}{\partial f_K} \right]^{-1} = \prod_{k=1}^K (f_{k-1} - F_{K+1}) = (1 - F_{K+1})^K \prod_{k=1}^K g_k^{K-k}.$$

We deduce that the density of  $(G_1, \dots, G_K) | \{\bar{t}_1, \dots, \bar{t}_K, F_{K+1}\}$  at point  $g$  with respect to Lebesgue measure is

$$\mathbf{p}(g) := (\text{const}) \mathbb{1}_{g \in (0,1)^K} \prod_{k=1}^K g_k^{K-k} \cosh [\mu_{\bar{t}_k}^0 ((1 - F_{K+1}) \prod_{\ell=1}^k g_\ell + F_{K+1})], \quad (39)$$

where we have used (36). From (18) and (33), one has

$$\hat{p}_k = 1 - \mathbf{F}_{\beta(1, K-k+1)} \left( \frac{F_k - F_{K+1}}{F_{k-1} - F_{K+1}} \right) = 1 - \mathbf{F}_{\beta(1, K-k+1)}(G_k) \quad (40)$$

where  $\mathbf{F}_\beta$  is the cumulative distribution of the Beta distribution in reference. We deduce that for any  $v \in (0, 1)$  and for any  $\ell \in [K]$ ,

$$\hat{p}_k = v \Leftrightarrow (G_1, \dots, G_K) \in [0, 1]^K \cap \{\mathbf{F}_{\beta(1, K-k+1)}^{-1}(1 - v) = G_k\},$$

so that

$$\bar{\mathbb{P}}[(\hat{p}_1, \dots, \hat{p}_K) \in \Gamma | \hat{p}_{k^0} \leq u] = \bar{\mathbb{P}}[(G_1, \dots, G_K) \in \bar{\Gamma} | G_{k^0} \geq \mathbf{F}_{\beta(1, K-\ell+1)}^{-1}(1 - u)], \quad (41)$$

where  $\bar{\Gamma}$  can be proved to be a **nonincreasing Borel set** from (40).

◦ **Step 2:** Let  $0 < x < y < 1$  and denote by  $\mu_x$  the following conditional law

$$\mu_x := \text{law}[(G_1, \dots, G_K) | \{\bar{t}_1, \dots, \bar{t}_K, F_{K+1}, G_{k^0} \geq x\}].$$

Remark that if there exists a measurable  $T : [0, 1]^K \mapsto [0, 1]^K$  such that

- $T$  is nondecreasing, meaning that for any  $g \in [0, 1]^K$ ,  $T(g) \geq g$ ;
- $T$  is such that push-forward of  $\mu_x$  by  $T$  gives  $\mu_y$ , namely  $T_{\#}\mu_x = \mu_y$ ;

then it holds

- $\mathbb{1}_{\{T(g) \in \bar{\Gamma}\}} \leq \mathbb{1}_{\{g \in \bar{\Gamma}\}}$ ;
- $\text{law}[T(G) | \{\bar{t}_1, \dots, \bar{t}_K, F_{K+1}, G_{k^0} \geq x\}] = \text{law}[G | \{\bar{t}_1, \dots, \bar{t}_K, F_{K+1}, G_{k^0} \geq y\}]$  where  $G = (G_1, \dots, G_K)$ .

In this case, we deduce that

$$\bar{\mathbb{P}}[G \in \bar{\Gamma} | G_{k^0} \geq x] \geq \bar{\mathbb{P}}[T(G) \in \bar{\Gamma} | G_{k^0} \geq x] = \bar{\mathbb{P}}[G \in \bar{\Gamma} | G_{k^0} \geq y].$$

If one can prove that such function  $T$  exists for any  $0 < x < y < 1$ , it proves that

$$x \mapsto \bar{\mathbb{P}}[G \in \bar{\Gamma} | G_{k^0} \geq x] \text{ is nonincreasing,}$$

and, in view of (41), it proves (38). Proving that such function  $T$  exists is done in the next step.

◦ **Step 3:** Let  $0 < x < y < 1$ . Consider the **Knothe-Rosenblatt transport map**  $T$  of  $\mu_x$  toward  $\mu_y$  following the order

$$k^0 \rightarrow k^0 + 1 \rightarrow \cdots \rightarrow K \rightarrow k^0 - 1 \rightarrow k^0 - 2 \rightarrow \cdots \rightarrow 1.$$

It is based on a sequence of conditional quantile transforms defined following the ordering above. Its construction is presented for instance in [24, Sec.2.3, P.67] or [31, P.20]. The transport  $T$  is defined as follows. Given  $z, z' \in [0, 1]^K$  such that  $z' = T(z)$  it holds

$$\begin{aligned} z'_{k^0} &= T^{(k^0)}(z_{k^0}); \\ z'_{k^0+1} &= T^{(k^0+1)}(z_{k^0+1}, z'_{k^0}); \\ &\vdots \\ z'_K &= T^{(K)}(z_K, z'_{K-1}, \dots, z'_{k^0}); \\ z'_{k^0-1} &= T^{(k^0-1)}(z_{k^0-1}, z'_K, \dots, z'_{k^0}); \\ &\vdots \\ z'_1 &= T^{(1)}(z_1, z'_2, \dots, z'_{k^0-1}, z'_K, \dots, z'_{k^0}); \end{aligned}$$

where  $T^{(k^0)}, T^{(k^0+1)}, \dots, T^{(K)}, T^{(k^0-1)}, \dots, T^{(1)}$  will be build in the sequel, in which we will drop their dependencies in the  $z'_k$ 's to ease notations. It remains to prove that

- $T$  is nondecreasing, meaning that for any  $g \in [0, 1]^K$ ,  $T(g) \geq g$ ;
- $T$  is such that push-forward of  $\mu_x$  by  $T$  gives  $\mu_y$ , namely  $T_{\#}\mu_x = \mu_y$ ;

to conclude. The last point is a property of the Knothe-Rosenblatt transport map. Proving the first point will be done in the rest of the proof.

◦ *Step 3.1:* We start by the first transport map  $T^{(k^0)} : [0, 1] \mapsto [0, 1]$ . Denote  $\mu_x^{(k^0)}$  the following conditional law

$$\mu_x^{(k^0)} := \text{law}[G_{k^0} | \{\bar{v}_1, \dots, \bar{v}_K, F_{K+1}, G_{k^0} \geq x\}],$$

and  $\mathbb{F}_x^{(k^0)}$  its cdf. Note that the Knothe-Rosenblatt construction gives  $T^{(k^0)} = (\mathbb{F}_y^{(k^0)})^{-1} \circ \mathbb{F}_x^{(k^0)}$ . We would like to prove that  $T^{(k^0)}(t) \geq t$  for all  $z \in (0, 1)$ . This is equivalent to prove that it holds  $\mathbb{F}_x^{(k^0)} \geq \mathbb{F}_y^{(k^0)}$ . For  $t \leq y$ ,  $\mathbb{F}_y^{(k^0)}(t) = 0$  and it implies that  $\mathbb{F}_x^{(k^0)}(t) \geq \mathbb{F}_y^{(k^0)}(t)$ . Let  $t > y$ , using the conditional density  $\mathbf{p}$  defined in (39), note that

$$\begin{aligned} \mathbb{F}_x^{(k^0)}(t) \geq \mathbb{F}_y^{(k^0)}(t) &\Leftrightarrow \frac{\int_x^t \mathbf{p}}{\int_x^1 \mathbf{p}} \geq \frac{\int_y^t \mathbf{p}}{\int_y^1 \mathbf{p}} \\ &\Leftrightarrow \int_x^t \int_y^1 \mathbf{p} \otimes \mathbf{p} \geq \int_y^t \int_x^1 \mathbf{p} \otimes \mathbf{p}, \end{aligned}$$

where, for example

$$\int_x^t \text{ means the integral over the hyper rectangle } [x, t] := \{(g_1, \dots, g_K) \in [0, 1]^K : x \leq g_{k^0} \leq t\}.$$

A simple calculation (see also Figure 6) gives that

$$\int_x^t \int_y^1 \mathbf{p} \otimes \mathbf{p} = \int_y^t \int_x^1 \mathbf{p} \otimes \mathbf{p} + \int_{[x, y] \times [t, 1]} \mathbf{p} \otimes \mathbf{p},$$

and it proves that  $\mathbb{F}_x^{(k^0)} \geq \mathbb{F}_y^{(k^0)}$ .

◦ *Step 3.2:* We continue with the second transport map in Knothe-Rosenblatt construction. Let  $z_{k^0} \in (x, 1)$  and denote  $\mu_{z_{k^0}}^{(k^0+1)}$  the following conditional law

$$\mu_{z_{k^0}}^{(k^0+1)} := \text{law}[G_{k^0+1} | \{\bar{v}_1, \dots, \bar{v}_K, F_{K+1}, G_{k^0} = z_{k^0}\}],$$

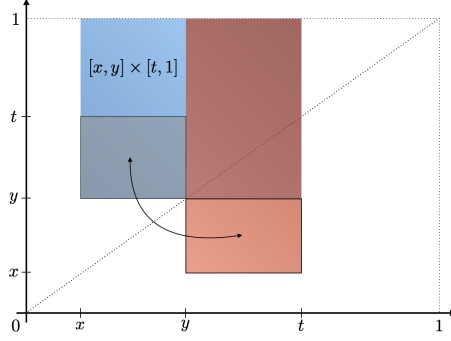


Figure 6. Note that, by symmetry the two boxed regions have same  $\mathbf{p} \otimes \mathbf{p}$  measure. The blue region is  $[x, t] \times [y, 1]$ , its measure is the measure of the red region (namely  $cD_y(t) \times D_x(1)$ ) more the bluest upper left corner (namely  $[x, y] \times [t, 1]$ ).

and  $\mathbb{F}_{z_{k^0}}^{(k^0+1)}$  its cdf. Let  $z'_{k^0} := T^{(k^0)}(z_{k^0})$  and denote  $\mu_{z'_{k^0}}^{(k^0+1)}$  the following conditional law

$$\mu_{z'_{k^0}}^{(k^0+1)} := \text{law}[G_{k^0+1} | \{\bar{v}_1, \dots, \bar{v}_K, F_{K+1}, G_{k^0} = z'_{k^0}\}],$$

and  $\mathbb{F}_{z'_{k^0}}^{(k^0+1)}$  its cdf. Note that  $x < z_{k^0} \leq z'_{k^0} = T^{(k^0)}(z_{k^0}) \leq 1$ . Again, we would like to prove that  $\mathbb{F}_{z_{k^0}}^{(k^0+1)} \geq \mathbb{F}_{z'_{k^0}}^{(k^0+1)}$  which implies that the transport map  $T^{(k^0+1)} := (\mathbb{F}_{z'_{k^0}}^{(k^0+1)})^{-1} \circ \mathbb{F}_{z_{k^0}}^{(k^0+1)}$  satisfies  $T^{(k^0+1)}(u) \geq u$  for all  $u \in (0, 1)$ .

Recall that the conditional density  $\mathbf{p}$  of  $G | \{\bar{v}_1, \dots, \bar{v}_K, F_{K+1}\}$  is given by (39) and recall that  $k^0 \in I_0$ . Observe that  $\mu_{k^0}^0 = 0$ , so that the conditional density of  $G | \{\bar{v}_1, \dots, \bar{v}_K, F_{K+1}, G_{k^0} = \mathbf{z}\}$  is

$$\begin{aligned} & (\text{const}) \mathbb{1}_{g \in (0,1)^K} \mathbb{1}_{g_{k^0} = \mathbf{z}} \prod_{k < k^0} g_k^{K-k} \cosh[\mu_{\bar{v}_k}^0 ((1 - F_{K+1}) \prod_{\ell=1}^k g_\ell + F_{K+1})] \\ & \times \prod_{k > k^0} g_k^{K-k} \cosh[\mu_{\bar{v}_k}^0 ((1 - F_{K+1}) \mathbf{z} \prod_{1 \leq \ell \neq k^0 \leq k} g_\ell + F_{K+1})]. \end{aligned} \quad (42)$$

Set  $\tau := z'_{k^0}/z_{k^0} \geq 1$  and  $G'_{k^0+1} = \tau G_{k^0+1}$  so that

$$z_{k^0} G_{k^0+1} = z'_{k^0} G'_{k^0+1}.$$

Denote  $G' := (G_1, \dots, G_{k^0-1}, G'_{k^0+1}, G_{k^0+2}, \dots, G_K) \in (0, 1)^{k^0-1} \times (0, \tau) \times (0, 1)^{K-k^0-1}$  and note that the conditional density of  $G' | \{\bar{v}_1, \dots, \bar{v}_K, F_{K+1}, G_{k^0} = \tau \mathbf{z}\}$  is

$$\begin{aligned} & (\text{const}) \mathbb{1}_{g \in (0,1)^{k^0-1} \times (0,\tau) \times (0,1)^{K-k^0-1}} \prod_{k < k^0} g_k^{K-k} \cosh[\mu_{\bar{v}_k}^0 ((1 - F_{K+1}) \prod_{\ell=1}^k g_\ell + F_{K+1})] \\ & \times \prod_{k > k^0} g_k^{K-k} \cosh[\mu_{\bar{v}_k}^0 ((1 - F_{K+1}) \mathbf{z} \prod_{1 \leq \ell \neq k^0 \leq k} g_\ell + F_{K+1})], \end{aligned}$$

which, up to some normalising constant, is the same as (42) up to the following change of support

$$\mathbb{1}_{g \in (0,1)^K} \leftrightarrow \mathbb{1}_{g' \in (0,1)^{k^0-1} \times (0,\tau) \times (0,1)^{K-k^0-1}}.$$

By an abuse of notation, we denote by  $\mathbf{p}$  this function, namely

$$\begin{aligned} \mathbf{p}(g) &= \prod_{k < k^0} g_k^{K-k} \cosh[\mu_{\bar{v}_k}^0 ((1 - F_{K+1}) \prod_{\ell=1}^k g_\ell + F_{K+1})] \\ & \times \prod_{k > k^0} g_k^{K-k} \cosh[\mu_{\bar{v}_k}^0 ((1 - F_{K+1}) \mathbf{z} \prod_{1 \leq \ell \neq k^0 \leq k} g_\ell + F_{K+1})]. \end{aligned}$$

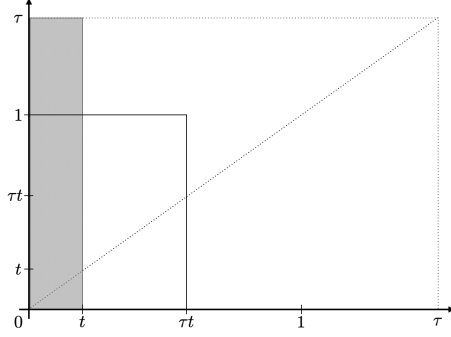


Figure 7. The two boxed rectangles have Lebesgue measure, namely  $\tau t$ . The  $\mathbf{p} \otimes \mathbf{p}$  measure of the grey box is greater than the  $\mathbf{p} \otimes \mathbf{p}$  measure of the white box.

We deduce that

$$\begin{aligned}
\mathbb{F}_{z_{k^0}}^{(k^0+1)}(t) \geq \mathbb{F}_{z'_{k^0}}^{(k^0+1)}(t) &\Leftrightarrow \overline{\mathbb{P}}(G_{k^0+1} \leq t | G_{k^0} = z_{k^0}) \geq \overline{\mathbb{P}}(G_{k^0+1} \leq t | G_{k^0} = z'_{k^0}) \\
&\Leftrightarrow \overline{\mathbb{P}}(G_{k^0+1} \leq t | G_{k^0} = z_{k^0}) \geq \overline{\mathbb{P}}(G'_{k^0+1} \leq \tau t | G_{k^0} = \tau z_{k^0}) \\
&\Leftrightarrow \frac{\int_{\mathcal{D}(t)} \mathbf{p}}{\int_{\mathcal{D}(1)} \mathbf{p}} \geq \frac{\int_{\mathcal{D}(\tau t)} \mathbf{p}}{\int_{\mathcal{D}(\tau)} \mathbf{p}} \\
&\Leftrightarrow \int_{\mathcal{D}(t) \times \mathcal{D}(\tau)} \mathbf{p} \otimes \mathbf{p} \geq \int_{\mathcal{D}(\tau t) \times \mathcal{D}(1)} \mathbf{p} \otimes \mathbf{p}, \tag{43}
\end{aligned}$$

where

$$\mathcal{D}(s) := \left\{ (g_1, \dots, g_{k^0-1}, g_{k^0+1}, \dots, g_K) \in (0, 1)^{K-1} : 0 < g_{k^0+1} \leq s \right\}.$$

We now present an inequality on the to conclude. Observe that we are integrating on domains depicted in Figure 7. The two boxes have same area for the uniform measure and we would like to compare their respective measure for the  $\mathbf{p} \otimes \mathbf{p}$  measure. We start by the next lemma whose proof is omitted.

**Lemma 13.** *Let  $a, b \geq 0$ . The function*

$$z \mapsto \cosh(az + b) \times \cosh(a/z + b)$$

*is non-decreasing on the domain  $[1, \infty)$ .*

Now, let  $(g_1, \dots, g_{k^0-1}, g_{k^0+2}, \dots, g_K) \in (0, 1)^{K-1}$  be fixed in the integrals (43). We are the looking at the weights of the domains  $(h_1, h_2) \in (0, t) \times (0, \tau)$  and  $(h_3, h_4) \in (0, \tau t) \times (0, 1)$  for the weight function  $w$  given by

$$\begin{aligned}
w(h_1, h_2) = & C_1 h_1^{K-k^0-1} \cosh \left[ \mu_{i_k}^0 ((1 - F_{K+1}) z_{k^0} \prod_{1 \leq \ell < k^0} g_\ell \times h_1 + F_{K+1}) \right] \\
& \times \prod_{k > k^0+1} \cosh \left[ \mu_{i_k}^0 ((1 - F_{K+1}) z_{k^0} \prod_{1 \leq \ell \neq k^0, k^0+1 \leq k} g_\ell \times h_1 + F_{K+1}) \right] \\
& \times h_2^{K-k^0-1} \cosh \left[ \mu_{i_k}^0 ((1 - F_{K+1}) z_{k^0} \prod_{1 \leq \ell < k^0} g_\ell \times h_2 + F_{K+1}) \right] \\
& \times \prod_{k > k^0+1} \cosh \left[ \mu_{i_k}^0 ((1 - F_{K+1}) z_{k^0} \prod_{1 \leq \ell \neq k^0, k^0+1 \leq k} g_\ell \times h_2 + F_{K+1}) \right],
\end{aligned}$$

where the constant  $C_1$  depends on  $(g_1, \dots, g_{k^0-1}, g_{k^0+2}, \dots, g_K) \in (0, 1)^{K-1}$ . By the change of variables  $h'_1 = h_3/t$  and  $h'_2 = th_4$ , the right hand term of (43) is given by the integration on the

domain  $(h'_1, h'_2) \in (0, t) \times (0, \tau)$  of the weight function  $w'$  given by

$$\begin{aligned}
w'(h'_1, h'_2) = & C_1 h'_1{}^{K-k^0-1} \cosh [\mu_{i_k}^0 ((1 - F_{K+1}) z_{k^0} \prod_{1 \leq \ell < k^0} g_\ell \times t \times h'_1 + F_{K+1})] \\
& \times \prod_{k > k^0+1} \cosh [\mu_{i_k}^0 ((1 - F_{K+1}) z_{k^0} \prod_{1 \leq \ell \neq k^0, k^0+1 \leq k} g_\ell \times t \times h'_1 + F_{K+1})] \\
& \times h'_2{}^{K-k^0-1} \cosh [\mu_{i_k}^0 ((1 - F_{K+1}) z_{k^0} \prod_{1 \leq \ell < k^0} g_\ell \times h'_2/t + F_{K+1})] \\
& \times \prod_{k > k^0+1} \cosh [\mu_{i_k}^0 ((1 - F_{K+1}) z_{k^0} \prod_{1 \leq \ell \neq k^0, k^0+1 \leq k} g_\ell \times h'_2/t + F_{K+1})].
\end{aligned}$$

Now, invoke Lemma 13 with

$$\begin{aligned}
a &= \mu_{i_k}^0 (1 - F_{K+1}) z_{k^0} \prod_{1 \leq \ell < k^0} g_\ell \times h \\
b &= \mu_{i_k}^0 F_{K+1} \\
z &= t \geq 1,
\end{aligned}$$

where  $h = h_1$  or  $h_2$ , to get that  $w' \geq w$  and so

$$\int_{\mathcal{D}(t) \times \mathcal{D}(\tau)} \mathbf{p} \otimes \mathbf{p} \geq \int_{\mathcal{D}(\tau t) \times \mathcal{D}(1)} \mathbf{p} \otimes \mathbf{p},$$

which concludes this part of the proof.

◦ *Step 3.3:* We continue by induction with the other transport maps in Knothe-Rosenblatt's construction. Assume that we have built  $z' := (z'_k, \dots, z'_{k^0})$  and  $z := (z_k, \dots, z_{k^0})$  for some  $k > k^0$ . Denote  $\mu_z^{(k+1)}$  the following conditional law

$$\mu_z^{(k+1)} := \text{law}[G_{k+1} | \{\bar{i}_1, \dots, \bar{i}_K, F_{K+1}, \underbrace{G_k = z_k, \dots, G_{k^0} = z_{k^0}}_{\text{denoted } G^{[k, k^0]} = z}\}],$$

and  $\mathbb{F}_z^{(k+1)}$  its cdf. Denote  $\mu_{z'}^{(k+1)}$  the following conditional law

$$\mu_{z'}^{(k+1)} := \text{law}[G_{k+1} | \{\bar{i}_1, \dots, \bar{i}_K, F_{K+1}, \underbrace{G_k = z'_k, \dots, G_{k^0} = z'_{k^0}}_{G^{[k, k^0]} = z'}\}],$$

and  $\mathbb{F}_{z'}^{(k+1)}$  its cdf. Note that  $z \leq z' = T^{(k)}(z) \leq 1$ . Again, we would prove that  $\mathbb{F}_z^{(k+1)} \geq \mathbb{F}_{z'}^{(k+1)}$  which implies that the transport map  $T^{(k+1)} := (\mathbb{F}_{z'}^{(k+1)})^{-1} \circ \mathbb{F}_z^{(k+1)}$  satisfies  $T^{(k+1)}(u) \geq u$  for all  $u \in (0, 1)$ .

For  $\mathbf{z} \in (0, 1)^{k-k^0} \times (x, 1)$ , the conditional density of  $G | \{\bar{i}_1, \dots, \bar{i}_K, F_{K+1}, G^{[k, k^0]} = \mathbf{z}\}$  is

$$\begin{aligned}
& (\text{const}) \mathbb{1}_{g \in (0, 1)^\kappa} \mathbb{1}_{g^{[k, k^0]} = \mathbf{z}} \prod_{m < k^0} g_m^{K-m} \cosh [\mu_{i_m}^0 ((1 - F_{K+1}) \prod_{\ell=1}^m g_\ell + F_{K+1})] \\
& \times \prod_{k^0 \leq m \leq k} z_m^{K-m} \cosh [\mu_{i_m}^0 ((1 - F_{K+1}) \prod_{1 \leq \ell < k^0} g_\ell \prod_{n=k^0}^m z_n + F_{K+1})] \\
& \times \prod_{k < m} g_m^{K-m} \cosh [\mu_{i_m}^0 ((1 - F_{K+1}) \prod_{1 \leq \ell < k^0} g_\ell \prod_{n=k^0}^k z_n \prod_{k < \ell \leq m} g_\ell + F_{K+1})].
\end{aligned}$$



Set  $\tau := \prod_{n=k^0}^k z'_n / \prod_{n=k^0}^k z_n \geq 1$  and  $G'_k = \tau G_{k^0+1}$  so that

$$\left[ \prod_{n=k^0}^k z'_n \right] G_{k+1} = \left[ \prod_{n=k^0}^k z_n \right] G'_{k+1}.$$

Then the proof follows the same idea as in *Step 3.2* and we will not detail it here.

◦ *Step 3.4:* This is the last step of the proof. Assume that we have built  $z' := (z'_K, \dots, z'_{k^0})$  and  $z := (z_K, \dots, z_{k^0})$ . Denote  $\mu_z^{(k^0-1)}$  the following conditional law

$$\mu_z^{(k^0-1)} := \text{law}[G_{k^0-1} | \{\bar{v}_1, \dots, \bar{v}_K, F_{K+1}, G^{[K, k^0]} = z\}],$$

and  $\mathbb{F}_z^{(k^0-1)}$  its cdf. Denote  $\mu_{z'}^{(k^0-1)}$  the following conditional law

$$\mu_{z'}^{(k^0-1)} := \text{law}[G_{k^0-1} | \{\bar{v}_1, \dots, \bar{v}_K, F_{K+1}, G^{[K, k^0]} = z'\}],$$

and  $\mathbb{F}_{z'}^{(k^0-1)}$  its cdf. Note that  $z \leq z' = T^{(K)}(z) \leq 1$ . Again, we would prove that  $\mathbb{F}_z^{(k^0-1)} \geq \mathbb{F}_{z'}^{(k^0-1)}$  which implies that the transport map  $T^{(k^0-1)} := (\mathbb{F}_{z'}^{(k^0-1)})^{-1} \circ \mathbb{F}_z^{(k^0-1)}$  satisfies  $T^{(k^0-1)}(u) \geq u$  for all  $u \in (0, 1)$ .

For  $\mathbf{z} \in (0, 1)^{K-k^0} \times (x, 1)$ , the conditional density of  $G | \{\bar{v}_1, \dots, \bar{v}_K, F_{K+1}, G^{[K, k^0]} = \mathbf{z}\}$  is

$$\begin{aligned} & (\text{const}) \mathbb{1}_{g \in (0, 1)^K} \mathbb{1}_{g^{[K, k^0]} = \mathbf{z}} \prod_{m < k^0} g_m^{K-m} \cosh [\mu_{\bar{v}_m}^0 ((1 - F_{K+1}) \prod_{\ell=1}^m g_\ell + F_{K+1})] \\ & \times \prod_{k^0 \leq m \leq K} \mathbf{z}_m^{K-m} \cosh [\mu_{\bar{v}_m}^0 ((1 - F_{K+1}) \prod_{1 \leq \ell < k^0} g_\ell \prod_{n=k^0}^m \mathbf{z}_n + F_{K+1})]. \end{aligned}$$

Now, let  $(g_1, \dots, g_{k^0-2}) \in (0, 1)^{k^0-2}$  be fixed and denote by

$$\begin{aligned} \forall g \in (0, 1), \quad w_z(g) &:= g^{K-k^0+1} \cosh [\mu_{\bar{v}_{k^0-1}}^0 ((1 - F_{K+1}) \prod_{\ell=1}^{k^0-2} g_\ell \times g + F_{K+1})] \\ & \times \prod_{k^0 \leq m \leq K} \mathbf{z}_m^{K-m} \cosh [\mu_{\bar{v}_m}^0 ((1 - F_{K+1}) \prod_{n=k^0}^m \mathbf{z}_n \prod_{\ell=1}^{k^0-2} g_\ell \times g + F_{K+1})]. \end{aligned}$$

and, substituting  $z$  by  $z'$ , define  $w_{z'}$  as well. Let  $t \in (0, 1)$ . Following the idea of *Step 3.2*, one can check that it is sufficient to prove that

$$\int_0^t \left( \int_0^1 w_z(g) w_{z'}(g') dg' \right) dg \geq \int_0^1 \left( \int_0^t w_z(g) w_{z'}(g') dg' \right) dg.$$

Substituting

$$\int_0^t \left( \int_0^t w_z(g) w_{z'}(g') dg' \right) dg$$

on both parts, one is reduced to prove that

$$\int_0^t \left( \int_t^1 w_z(g) w_{z'}(g') dg' \right) dg \geq \int_0^t \left( \int_t^1 w_{z'}(g) w_z(g') dg' \right) dg.$$

Observe that  $g \leq g'$  in the last two integrals. Now, we have this lemma whose proof is omitted.

**Lemma 14.** Let  $0 < a \leq a'$  and  $b > 0$ . The function

$$z \mapsto \frac{\cosh(az + b)}{\cosh(a'z + b)}$$

is non-increasing on the domain  $(0, \infty)$ .

Let  $g \leq g'$ . From Lemma 14, we deduce that  $\cosh(ag + b) \cosh(a'g' + b) \geq \cosh(ag' + b) \cosh(a'g + b)$ , proving that  $w_z(g)w_{z'}(g') \geq w_{z'}(g)w_z(g')$ . It proves that  $T^{(k^0-1)}(u) \geq u$  for all  $u \in (0, 1)$ .

We then proceed by induction for  $k^0 - 1 \rightarrow k^0 - 2 \rightarrow \dots \rightarrow 1$ . The proof follows the same line as above, Step 3.4.

sectionProof of Theorem 5 From Corollary 4 we deduce that, under  $\mathbb{H}_0$  and conditionally on the selection event  $\{\hat{m} = a\}$  with  $a \leq K - 1$ , (17) remains true, giving the conditional cumulative distribution function of  $\lambda_b$ . As a consequence and under the same conditioning.

$$\frac{F_{abc}(\lambda_b)}{F_{abc}(\lambda_a)} \sim \mathcal{U}(0, 1).$$

Since the conditional distribution doesn't depend on the condition, it is also the unconditional distribution. Finally considerations of distribution under the alternative show that to obtain a  $p$ -value we must consider the complement to 1 of the quantity above.

## Appendix D: Main notation

$[a]$	the set of integers $\{1, \dots, a\}$ ,
$X$	a full rank $n \times p$ design matrix,
$\sigma^2$	the variance of the errors,
$\bar{i}_k; \varepsilon_k$	the indexes and the signs of the variables that enter in the LAR path,
$\hat{i}_k$	a way of coding both indexes and signs, see (1),
$\bar{S}^k$	$\{\bar{i}_1, \dots, \bar{i}_k\}$ , a possible selected support,
$S_0$	the true support,
$E_k$	$\text{Span}(X_{\bar{i}_1}, \dots, X_{\bar{i}_k})$ ,
$P_k(P_k^\perp)$	projector on (the orthogonal of) $E_k$ ,
$\gamma_{m_k, v_k^2}(\phi_{m_k, v_k^2})$	The normal distribution (density) with mean $m_k$ given (13) and variance $v_k^2 = \sigma^2 \rho_k^2$ given by (15),
$Z$	the vector of correlations, obtained by sym. from $\bar{Z}$ defined by (9),
$R$	the variance-covariance matrix of $Z$ , see (10),
$K$	the number of knots that are considered, see Figure 2, and Section 2.5
$K_{\text{select}}$	see Figure 2
$\hat{m}$	the chosen size ,
$\hat{S}$	the chosen set of variables : $\bar{S}^{\hat{m}}$ ,
$\rho_k^2$	$v_k^2 / \sigma^2$ , see (13),
$\theta_j(i_1, \dots, i_k)$	see (11),
$\theta^\ell$	$\theta(\hat{i}_1, \dots, \hat{i}_\ell)$ ,
$M_{i_1, \dots, i_\ell}$	the submatrix of $R$ obtained by keeping the columns, and the rows indexed by $\{i_1, \dots, i_\ell\}$ ,
$F_i; \mathcal{P}_{ij}$	$F_i := \Phi_i(\lambda_i) := \Phi(\lambda_i / (\sigma \rho_i))$ and $\mathcal{P}_{ij}$ is given by (14),
$\mathbb{F}_{abc}(t)$	see (16)
$\hat{\alpha}_{abc}$	the $p$ -value of the <i>generalized spacing test</i> , see (18),
$\mathcal{S}_{abc}$	$\mathbb{1}_{\{\hat{\alpha}_{abc} \leq \alpha\}}$ , the <i>generalized spacing test</i> see (19).

## Appendix E: Cubature of integral by lattice rule

Our goal is to compute the integral of some function  $f$  on the hypercube of dimension  $d$ , namely

$$I := \int_{[0,1]^d} f(x) dx.$$

We want to approximate it by a finite sum over  $n$  points

$$I_n := \frac{1}{n} \sum_{i=1}^n f(x^{(i)}).$$

A convenient way of constructing the sequence  $x^{(i)}, i = 1, \dots, n$  is the so-called *lattice rule*: from the first point  $x^{(1)}$  we deduce the others  $x^{(i)}$  by

$$x^{(i)} = \{i.x^{(1)}\},$$

where the  $\{\}$  brackets mean that we take the fractional part coordinate by coordinate. In such a case the error given by

$$E(f, n, x^{(1)}) = I - I_n$$

is a function, in particular, of starting point  $x^{(1)}$ .

The Fast-rank algorithm [21] is a fast algorithm that finds, component by component and as a function of the prime  $n$ , the sequence of coordinates of  $x^{(1)}$  that minimizes the maximal error when  $f$  varies in a unit ball  $\mathcal{E}$  of some *RKHS*, namely a tensorial product of *Koborov spaces*. In addition it gives an expression of its minimax error, namely

$$\max_{f \in \mathcal{E}} (f, n, x^{(1)}).$$

In practice, very few properties are known on the function  $f$ , so the result above is not directly applicable. Nevertheless for many functions  $f$ , it happens that the convergence of  $I_n$  to  $I$  is “fast”: typically of the order  $1/n$  while the Monte-Carlo method (choosing the  $x^{(i)}$  at random) converges at rate  $1/\sqrt{n}$ .

A reliable estimation of the error is obtained by adding a *Monte-Carlo layer* as in [16] for instance. This can be done as follows. Let  $U$  a **unique** uniform variable on  $[0, 1]^d$ , we define

$$x_U^{(i)} := \{i.x^{(1)} + U\}, \quad I_{n,U} := \frac{1}{n} \sum_{i=1}^n f(x_U^{(i)}).$$

Classical computations show that  $I_{n,U}$  is now an unbiased estimator of  $I$ . In a final step, we perform  $N$  (in practice 15-20) independent repetitions of the experiment above and we compute usual asymptotic confidence intervals for independent observations.