



HAL
open science

Trade-off between prediction and FDR for high-dimensional Gaussian model selection

Perrine Lacroix, Marie-Laure Martin

► **To cite this version:**

Perrine Lacroix, Marie-Laure Martin. Trade-off between prediction and FDR for high-dimensional Gaussian model selection. 2023. hal-03978309v1

HAL Id: hal-03978309

<https://hal.science/hal-03978309v1>

Preprint submitted on 8 Feb 2023 (v1), last revised 11 Apr 2024 (v3)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

TRADE-OFF BETWEEN PREDICTION AND FDR FOR HIGH-DIMENSIONAL GAUSSIAN MODEL SELECTION

Perrine Lacroix

Laboratoire de Mathématiques d'Orsay, CNRS, Université Paris-Saclay, Orsay, France
 Université Paris-Saclay, CNRS, INRAE, Université Evry, Institute of Plant Sciences Paris-Saclay (IPS2),
 91190, Gif sur Yvette, France
 Université Paris Cité, Institute of Plant Sciences Paris-Saclay (IPS2), 91190, Gif sur Yvette, France
perrine.lacroix@universite-paris-saclay.fr

Marie-Laure Martin

Université Paris-Saclay, AgroParisTech, INRAE, UMR MIA Paris-Saclay, 91120, Palaiseau, France
 Université Paris-Saclay, CNRS, INRAE, Université Evry, Institute of Plant Sciences Paris-Saclay (IPS2),
 91190, Gif sur Yvette, France
 Université Paris Cité, Institute of Plant Sciences Paris-Saclay (IPS2), 91190, Gif sur Yvette, France
marie-laure.martin@inrae.fr

ABSTRACT

In the context of the high-dimensional Gaussian linear regression for ordered variables, we study the variable selection procedure via the minimization of the penalized least-squares criterion. We focus on model selection where we propose to control predictive risk and False Discovery Rate simultaneously. For this purpose, we obtain a convenient trade-off thanks to a proper calibration of the hyperparameter K appearing in the penalty function. We obtain non-asymptotic theoretical bounds on the False Discovery Rate with respect to K . We then provide an algorithm for the calibration of K . It is based on completely observable quantities in view of applications. Our algorithm is validated by an extensive simulation study.

Keywords Ordered variable selection · Prediction · FDR · High-dimension · Gaussian regression · Hyperparameter calibration

1 Introduction.

1.1 The issue.

We consider the following high-dimensional Gaussian linear regression model:

$$Y = X\beta^* + \varepsilon. \quad (1.1)$$

The random response vector $Y = \left((Y_i)_{\{1 \leq i \leq n\}} \right)^T \in \mathbb{R}^n$ is regressed on p deterministic vectors: $X_1 = \left((x_{i1})_{\{1 \leq i \leq n\}} \right)^T, \dots, X_p = \left((x_{ip})_{\{1 \leq i \leq n\}} \right)^T$. The design matrix of size $n \times p$ is denoted by $X = (X_1, \dots, X_p)$. The noise $\varepsilon = \left((\varepsilon_i)_{\{1 \leq i \leq n\}} \right)^T$ is assumed to be Gaussian: $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ with $\sigma^2 > 0$. In the high-dimensional context, additional assumptions of regularity are required and we assume that β^* is sparse, meaning that only a few coefficients are non-zero. In the following, a variable X_j corresponding to a non-zero coefficient β_j^* is called an active variable. Otherwise the variable is said to be non active.

In this paper, we are interested in variable selection procedures. To the best of our knowledge, some procedures focus on the prediction of the response variable Y through a control of the predictive risk. Others focus on limiting the number of

selected non active variables through a control of the False Discovery Rate. There also exists procedures where several cost functions are considered simultaneously. In the line of the latter, our goal is to identify a set of variables from a model selection procedure by limiting the selection of non active variables while maintaining accurately prediction performances.

1.2 Related works.

In a variable selection procedure, a cost function has to be defined. The predictive risk (PR) and the False Discovery Rate (FDR) are the common used cost functions.

The penalized methods to control the predictive risk. The penalization procedure balances goodness of fit and sparsity: the smaller the penalty function, the better the fitting to the data but the higher the number of selected variables. In high-dimension, the most popular method is the Lasso criterion [1] where the estimator $\hat{\beta}_\lambda$ of β^* is the solution of:

$$\hat{\beta}_\lambda = \arg \min_{\beta \in \mathbb{R}^p} \left\{ \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\}, \quad (1.2)$$

where $\|\cdot\|_1$ and $\|\cdot\|_2$ design the ℓ_1 -norm and the euclidean norm of a vector respectively. The main challenge is to calibrate the hyperparameter $\lambda > 0$. If λ is chosen to be proportional to $\sigma \sqrt{\frac{\log(p)}{n}}$, then the predictive risk is bounded [2, 3]. However, the noise being usually unknown, the choice of λ remains tricky. Therefore, an alternative is to solve the Lasso criterion for a λ within a reasonable interval by using subsamples [4] or resamples [5]. The selected variables are then defined as the variables with the highest selection frequencies. Such alternative is no longer sensitive to the choice of λ but the main challenge lies in the threshold on the frequency defining the selected variables.

In this paper, we consider a model selection procedure composed of three steps. The first step consists in solving the Lasso criterion on a relevant grid Λ . Each $\lambda \in \Lambda$ defines a variable subset to get a collection \mathcal{M} of relevant subsets of variables with a wide range of sizes. In the second step, the least-squares estimator onto each variable subset of \mathcal{M} is calculated leading to a collection of estimators $(\hat{\beta}_m)_{m \in \mathcal{M}}$. Lastly, the following penalized least-squares minimization is solved to select the best m of \mathcal{M} :

$$\hat{m} = \arg \min_{m \in \mathcal{M}} \left\{ \|Y - X\hat{\beta}_m\|_2^2 + \text{pen}(D_m) \right\}, \quad (1.3)$$

where D_m is the dimension of the model m and the function pen is a penalty function increasing with D_m .

Selecting \hat{m} from \mathcal{M} by minimizing (1.3) corresponds to select $\hat{\lambda}$ from Λ by minimizing (1.2). Hence, the main challenge is the definition of pen that makes the best trade-off between goodness of fit and sparsity within \mathcal{M} . Among the most famous methods for model selection, we can cite V -fold cross-validation [6, 7], AIC [8], Cp-Mallows [9], BIC [10] and eBIC [11]. For these penalty functions, the predictive risk is bounded when σ^2 is known and when the sample size n tends to infinity. When n is fixed, relatively small, and possibly smaller than the dimension p , a non-asymptotic point of view is preferable to get properties for all couples of (n, p) . In this direction, [12] propose some penalty functions depending on the collection complexity such that \hat{m} guarantees non-asymptotic optimal control of the predictive risk. If the model collection is nested with a known variance, $\text{pen}(D_m) = 2\sigma^2 D_m$ allows to achieve an optimal non-asymptotic control of the predictive risk [8]. If the model collection is fixed and large (for instance with an exponential growth with D_m) and if the variance is unknown, this optimal control is obtained with the data-driven penalties [12, 13]. Lastly, if the model collection is data-dependent and if the variance σ^2 is unknown, the LinSelect penalty [14, 15] guarantees an optimal control of the predictive risk.

The multiple testing methods to control the False Discovery Rate. In the multiple testing procedure, the p tests $H_0 = \{\beta_j^* = 0\}$ versus $H_1 = \{\beta_j^* \neq 0\}$ are performed independently to get a list of p -values. Variables associated with a p -value smaller than a threshold are selected and the challenge is to find this threshold to obtain an upper bound on a function of the number of selected non active variables. First methods control the Family-Wise Error (FWER) which is the probability of selecting at least one non active variable [16, 17]. However, these methods tend to be conservative leading to a tiny set of selected variables. An alternative consists in controlling the FDR which is the expectation of the proportion of non active variables among the selected ones. The authors of [18] first provide a threshold assuming independence of the p -values. This hypothesis is then relaxed in [19, 20, 21, 22].

Instead of considering the p -values, the knockoff filter method [23] consists in introducing copies of X_j built to be non active variables to calibrate a threshold on test statistics.

The simultaneous control of several cost functions. Controlling PR or FDR is usually performed independently in the literature and yield different sets of selected variables. For a PR control, selected variables aim at correctly predict

a new observation of Y , without guarantying that the set of selected variables does not contain non active variables. Conversely, when the cost function is the FDR, the number of non active variables is controlled at the prize that some active variables are not selected.

Therefore, recent works have been proposed to combine prediction and FDR approaches to select all but only active variables. For instance, [24] propose a multi-step algorithm where a threshold procedure is applied on some Lasso estimators computed for specific values of λ . In addition to prediction performances, a consistency property on the selected variable set is satisfied under some conditions on X . Another idea is the post-selection inference [25, 26] where the principle is to test the relevance of each selected variable by a model selection procedure. Valid confidence intervals are provided from conditional hypothesis tests for each model of the collection in addition to a PR control. Their work has been generalized by [27, 28, 29] and a review can be found in [30].

In a completely different direction, [31, 32] propose to control the False Negative Rate (FNR) in addition to the FDR. A good FNR control ensures that most of the active variables are selected. So, minimizing a weighted sum of FDR and FNR provide a set of variables close to the set of active variables. However, improving FDR control deteriorates FNR control and vice versa. Hence, optimal controls of both criteria are impossible to achieve.

Some other papers propose to combine the FDR with the PR. An additive motivation to consider the PR is its behavior between the learning phase and the over-fitting phase. In the learning phase, the addition of a variable in the selected set drastically reduces the PR, whereas in the over-fitting phase, it increases proportionally to the noise level. Firstly, in the standard multivariate normal mean problem with a known variance, [33] propose a penalty function in the model selection procedure built from a multiple testing procedure. They obtain simultaneously sharp asymptotic minimality of the FDR and the PR. Then, [34] propose the Sorted ℓ_1 penalized estimator (SLOPE) which is the minimizer of the Lasso criterion (1.2) where λ is replaced by a p -vector built from a multiple testing procedure. When the variables X_1, \dots, X_p are orthogonal, they obtain a non-asymptotic control of the FDR additionally to the asymptotic minimax convergence rate of the PR. This asymptotic convergence of the FDR has been generalized under a wide range of hypotheses, for instance, for a random design in [35].

1.3 Main contributions.

The originality of this paper is to obtain a control of the FDR in addition to the PR control in model selection through a convenient calibration of the penalty.

We assume variables are ordered: the most relevant one to explain Y is X_1 ; the most relevant couple of two variables is (X_1, X_2) ; and so on. A natural model collection is the one containing the nested models respecting the variable order. This framework sounds restrictive but allows to derive theoretical expressions of the FDR in the considered model selection procedure. According to [36], all the penalty functions defined by:

$$\text{pen}(D_m) = K\sigma^2 D_m, \quad \forall m \in \mathcal{M}, \quad (1.4)$$

provide a non-asymptotic control of the PR for $K > 1$ when variables are ordered.

Theoretical bounds on the FDR in model selection: Although the model selection procedure is built for a PR control, we obtain non-asymptotic lower and upper bounds on the FDR with respect to $K > 0$ when σ^2 is known. We show that these bounds only involve some evaluations of cumulative functions of the standard Gaussian and of some chi-squared variables. Whatever the noise level, FDR is always strictly positive. When K tends to infinity, the FDR converges to 0 with an exponential rate. So, a low value of the FDR is satisfying as soon as the value of K is not too large.

Calibration of the hyperparameter K : The obtained theoretical bounds depend on the parameters β^* and σ^2 . We replace them with estimators to obtain completely data-dependent bounds on the FDR. Then, we propose a procedure to calibrate the hyperparameter K . Our algorithm is validated on an extensive simulation study.

1.4 Outline of the paper.

The rest of the paper is organized as follows. Section 2 introduces the Gaussian linear regression model and some notations. Section 3 contains theoretical results. Since an increase of the hyperparameter K leads to a decrease of the FDR, it motivates the study of the FDR function in model selection with respect to K . As the FDR has an intractable expression, bounds are obtained when variables are ordered and the variance is known. We establish an exponential convergence rate of the FDR function when K tends to infinity. The special case of orthogonal design matrix is studied to illustrate the main results. In Section 4, an algorithm is proposed to calibrate the hyperparameter K in the penalty function to get a convenient trade-off between FDR and PR controls. It is based on the simultaneous study of the estimated FDR upper bound and the estimated PR which depend on properly chosen estimators of σ^2 and β^* . Section 5 contains conclusions and perspectives. In Section 6, proofs of all the theoretical results are provided. Lastly, a validation

of the chosen estimators of σ^2 and β^* and of our algorithm to calibrate K is proposed in Section 7 through an extensive simulation study with different parameters.

2 Model and notations.

Let us consider the Gaussian linear regression model given in (1.1). We define $q = \min(n, p)$ and assume that (X_1, \dots, X_q) is a family of linearly independent vectors. We consider the deterministic and nested model collection of linear spaces:

$$\mathcal{M} = \left\{ m_0 = \{0\}, m_1 = \text{Span}(X_1), \dots, m_q = \text{Span}(X_1, X_2, \dots, X_q) \right\}. \quad (2.1)$$

By construction, the true model $m^* = \text{Span}(X_j, j \text{ s.t. } \beta_j^* \neq 0)$ belongs to \mathcal{M} .

For each $m \in \mathcal{M}$, D_m is the dimension of m and $\hat{\beta}_m$ is the least-squares estimator onto m :

$$\hat{\beta}_m = \arg \min_{\{\beta, X\beta \in m\}} \left\{ \|Y - X\beta\|_2^2 \right\}.$$

With the definition of q and properties on the family (X_1, \dots, X_q) , $\hat{\beta}_m$ is unique for each $m \in \mathcal{M}$.

For all $K > 0$, we define the function crit_K on \mathcal{M} as:

$$\text{crit}_K(m) = \|Y - X\hat{\beta}_m\|_2^2 + K\sigma^2 D_m,$$

and the selected model $\hat{m}(K)$ by:

$$\hat{m}(K) = \arg \min_{m \in \mathcal{M}} \left\{ \text{crit}_K(m) \right\}. \quad (2.2)$$

We define $\text{PR}(m)$ the predictive risk associated to the model $m \in \mathcal{M}$ by:

$$\text{PR}(m) = \mathbb{E} \left[\|Y - X\hat{\beta}_m\|_2^2 \right], \quad (2.3)$$

where \mathbb{E} designs the expectation under the distribution of Y satisfying (1.1). We define successively $FP(m)$ the number of variables contained in m but not in m^* , the false discovery proportion by:

$$\text{FDP}(m) = \frac{FP(m)}{\max(D_m, 1)};$$

and the False Discovery Rate by:

$$\text{FDR}(m) = \mathbb{E} \left[\text{FDP}(m) \right].$$

Finally, $\langle \cdot, \cdot \rangle$ designs the canonical scalar product in \mathbb{R}^n , $\Pi_{\mathcal{X}}$ denotes the orthogonal projection function onto the space \mathcal{X} , Φ designs the standard Gaussian cumulative distribution function and $F_{\chi^2(k)}$ is the cumulative distribution function of a chi-squared variable with k degrees of freedom. By convention, an intersection or an union from indices k to ℓ with $k > \ell$ are the intersection or the union over an empty set. In the same way, the set $\{k, \dots, \ell\}$ is empty if $k > \ell$.

3 The main results.

In this section, the variance σ^2 is supposed to be known. We first present intuitions that lead to study $\text{FDR}(\hat{m}(K))$ in model selection. Non-asymptotic bounds on $\text{FDR}(\hat{m}(K))$ are obtained in Theorem 3.2, as well as asymptotic behaviors when K tends to infinity in Corollary 3.4. Finally, the particular case where X is the orthogonal design matrix is studied to illustrate the main results.

3.1 Intuitions.

According to [12], the penalty function (1.4) satisfies a non-asymptotic control of the PR if and only if $K > 1$. The constant $K = 2$ allows to achieve the optimal asymptotic control of the PR. Hence, 2 is commonly chosen in practice but other values of K close to 2 can give equivalent even better non-asymptotic prediction performances. In this direction, we propose to calibrate the hyperparameter K among those leading to prediction performances close to or better than for $K = 2$ while satisfying a control of the FDR. The calibration is based on both $\text{PR}(\hat{m}(K))$ and $\text{FDR}(\hat{m}(K))$ functions with respect to K .

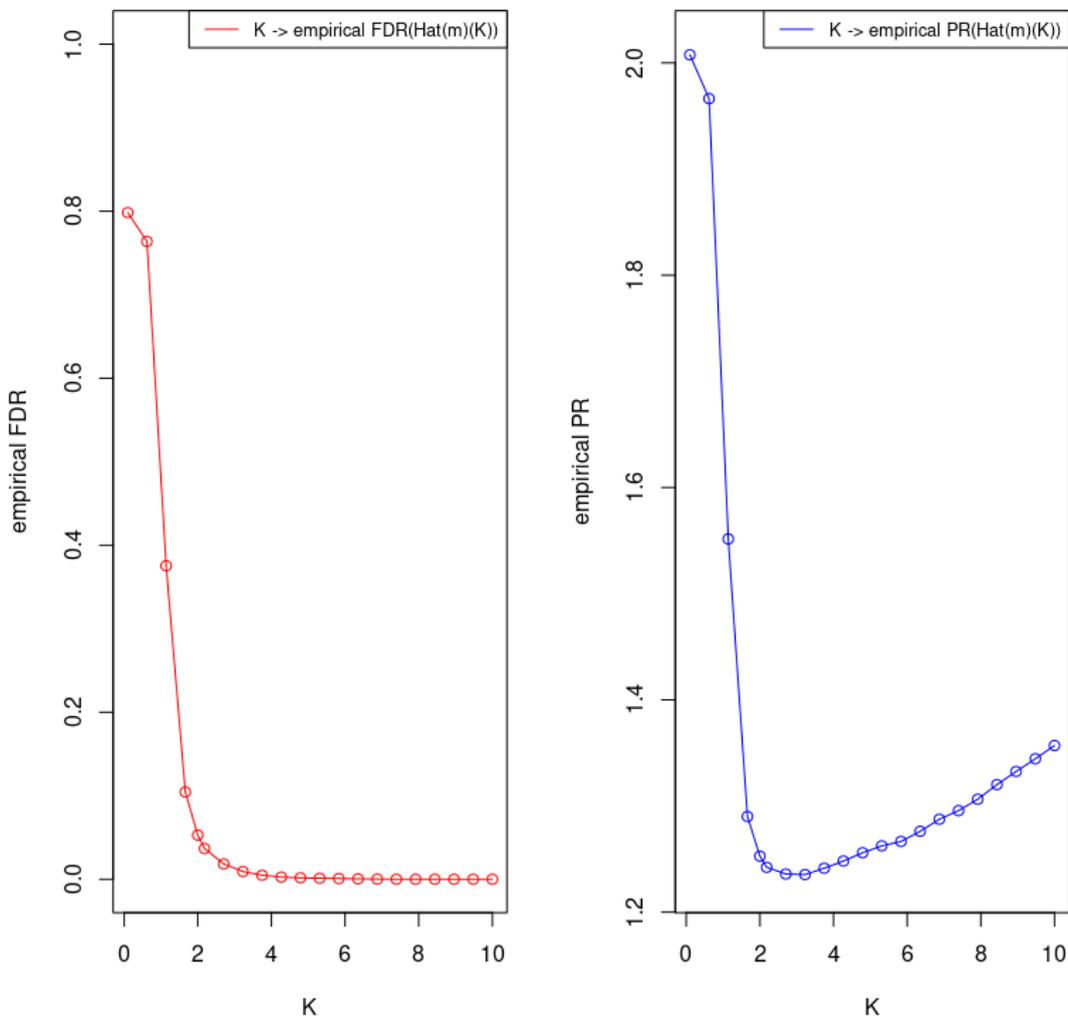


Figure 1: Curves of the empirical estimations of $\text{FDR}(\hat{m}(K))$ and $\text{PR}(\hat{m}(K))$ for the toy data set described in Subsection 7.1.

In Figure 1, we propose an illustration of our intuitions by plotting the empirical estimators of $\text{PR}(\hat{m}(K))$ and $\text{FDR}(\hat{m}(K))$ on a regular grid of positive K . Graphs are obtained from the *toy data set* described in Section 7. We observe that for all $K \in [2, 3]$, the empirical $\text{PR}(\hat{m}(K))$ values are kept low while the $\text{FDR}(\hat{m}(K))$ function decreases with K . Hence, in this example, the choice $K = 3$ is more judicious than $K = 2$ since it ensures a stronger control of the FDR while satisfying similar prediction performances.

Increasing the constant K to limit the non active variable selection is known for the asymptotic point of view. Indeed, AIC and Cp-Mallows penalties [8, 9], where K is fixed to 2, give asymptotically the best set of variables for prediction performances; while BIC penalty [10], where K is fixed to $\log(n)$, exactly recovers asymptotically the set of active variables. Obtaining the asymptotic properties of AIC, Cp-Mallows and BIC penalties simultaneously is impossible [37], but it suggests that a value of $K \in [2, \log(n)]$ would get reasonable (but not necessarily optimal) values for both PR and FDR in a non-asymptotic framework. In this way, we propose to study the function $\text{FDR}(\hat{m}(K))$ in the model selection procedure (1.4) when variables are ordered.

3.2 Bounds on the FDR in model selection.

3.2.1 FDR expression in model selection for ordered variables.

Let us assume that $K > 0$ and crit_K is injective on \mathcal{M} . If $D_m^* = q$, $\text{FDR}(\hat{m}(K)) = 0$. Otherwise, the $\text{FDR}(\hat{m}(K))$ is expressed within the model selection procedure as:

$$\text{FDR}(\hat{m}(K)) = \sum_{r=D_m^*+1}^q \frac{r-D_m^*}{r} \mathbb{P} \left(\left\{ \bigcap_{\substack{\ell=0 \\ \ell \neq r}}^q \{\text{crit}_K(m_r) < \text{crit}_K(m_\ell)\} \right\} \right). \quad (3.1)$$

A detailed proof of (3.1) can be found in Subsection 6.1.

By using the decomposition

$$\left\{ \bigcap_{\ell=0}^{r-1} \{\text{crit}_K(m_r) < \text{crit}_K(m_\ell)\} \right\} \cap \left\{ \bigcap_{\ell=r+1}^q \{\text{crit}_K(m_r) < \text{crit}_K(m_\ell)\} \right\},$$

we obtain the following proposition:

Proposition 3.1. *Let us consider the ordered variable framework and the model collection (2.1) where $q = \min(n, p)$, $m^* \in \mathcal{M}$ and $D_m^* < q$. Let us assume that crit_K is injective on \mathcal{M} . Let us apply the Gram–Schmidt process to obtain (u_1, \dots, u_q) the orthonormal basis of \mathbb{R}^q such that $\text{Span}(X_1, \dots, X_j) = \text{Span}(u_1, \dots, u_j)$, $\forall j \in \{1, \dots, q\}$. If $p < n$, (u_1, \dots, u_q) is naturally completed to an orthonormal basis (u_1, \dots, u_n) on \mathbb{R}^n by the incomplete basis theorem.*

Then, $\forall K > 0$,

$$\text{FDR}(\hat{m}(K)) = \sum_{r=D_m^*+1}^q \frac{r-D_m^*}{r} P_r(K) Q_r(K, \beta^*, \sigma^2), \quad (3.2)$$

where for each $r \in \{D_m^* + 1, \dots, q\}$,

$$P_r(K) = \mathbb{P} \left(\bigcap_{\ell=r+1}^q \left\{ \sum_{k=r+1}^{\ell} Z_k^2 < K(\ell - r) \right\} \right), \quad (3.3)$$

where $Z_k \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$, $\forall k \in \{r+1, \dots, q\}$,

$$\text{and } Q_r(K, \beta^*, \sigma^2) = \mathbb{P} \left(\bigcap_{\ell=0}^{r-1} \left\{ \sum_{k=\ell+1}^r \langle Y, u_k \rangle^2 > K\sigma^2(r - \ell) \right\} \right).$$

Proof of Proposition 3.1 can be found in Subsection 6.1.

3.2.2 General bounds.

In (3.2), the $P_r(K)$ terms do not depend on data. Conversely, the $Q_r(K, \beta^*, \sigma^2)$ terms depend on the data. Thus, to understand the behavior of the FDR function with respect to $\hat{m}(K)$, we propose to bound the $Q_r(K, \beta^*, \sigma^2)$ terms in the following theorem:

Theorem 3.2. *Let us consider the ordered variable framework and the model collection (2.1) where $q = \min(n, p)$. Let us suppose that $m^* \in \mathcal{M}$ and $D_m^* < q$. The notation Φ stands for the standard gaussian cumulative distribution function and $F_{\chi^2(k)}$ is the cumulative distribution function of a chi-squared variable with k degrees of freedom. Let us assume that $\forall K > 0$, crit_K is injective on \mathcal{M} . Let us apply the Gram–Schmidt process to obtain (u_1, \dots, u_q) the orthonormal basis of \mathbb{R}^q such that $\text{Span}(X_1, \dots, X_j) = \text{Span}(u_1, \dots, u_j)$, $\forall j \in \{1, \dots, q\}$. If $p < n$, (u_1, \dots, u_q) is naturally completed to an orthonormal basis (u_1, \dots, u_n) on \mathbb{R}^n by the incomplete basis theorem. Then, $\forall K > 0$, $\hat{m}(K)$ satisfies:*

$$b(K, \beta^*, \sigma^2) \leq \text{FDR}(\hat{m}(K)) \leq B(K, \beta^*, \sigma^2), \quad (3.4)$$

Trade-off between prediction and FDR for variable selection

where $[K \mapsto b(K, \beta^*, \sigma^2)]$ and $[K \mapsto B(K, \beta^*, \sigma^2)]$ are two real-valued functions on \mathbb{R}^+ defined by:

$$\begin{aligned} b(K, \beta^*, \sigma^2) &= \sum_{r=D_{m^*}+1}^q \left(\frac{r-D_{m^*}}{r} P_r(K) \underline{f}_r(K, \beta^*, \sigma^2) \right), \\ B(K, \beta^*, \sigma^2) &= \sum_{r=D_{m^*}+1}^q \left(\frac{r-D_{m^*}}{r} P_r(K) \bar{f}_r(K, \beta^*, \sigma^2) \right), \end{aligned} \quad (3.5)$$

with for all $K > 0$,

1. for each $r \in \{D_{m^*} + 1, \dots, q\}$,

$$P_r(K) = \mathbb{P} \left(\bigcap_{\ell=r+1}^q \left\{ \sum_{k=r+1}^{\ell} Z_k^2 < K(\ell-r) \right\} \right),$$

where $Z_k \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$, $\forall k \in \{r+1, \dots, q\}$.

2. for each $r \in \{D_{m^*} + 1, \dots, q\}$ and for all $\ell \in \{1, \dots, r\}$, $\underline{f}_\ell(\cdot, \beta^*, \sigma^2)$ is defined by:

$$\begin{aligned} \underline{f}_1(K, \beta^*, \sigma^2) &= G_1 \\ \underline{f}_\ell(K, \beta^*, \sigma^2) &= G_\ell + H_\ell \underline{f}_{\ell-1}(K, \beta^*, \sigma^2), \quad \forall \ell \in \{2, \dots, r\}, \end{aligned}$$

with for $\ell \in \{1, \dots, D_{m^*}\}$:

$$G_\ell = 2 - \left(\Phi \left(\sqrt{\ell K} - \frac{\langle X\beta^*, u_\ell \rangle}{\sigma} \right) + \Phi \left(\sqrt{\ell K} + \frac{\langle X\beta^*, u_\ell \rangle}{\sigma} \right) \right),$$

for $\ell \in \{2, \dots, D_{m^*}\}$:

$$\begin{aligned} H_\ell &= \Phi \left(\sqrt{\ell K} - \frac{\langle X\beta^*, u_\ell \rangle}{\sigma} \right) + \Phi \left(\sqrt{\ell K} + \frac{\langle X\beta^*, u_\ell \rangle}{\sigma} \right) \\ &\quad - \left(\Phi \left(\sqrt{K} - \frac{\langle X\beta^*, u_\ell \rangle}{\sigma} \right) + \Phi \left(\sqrt{K} + \frac{\langle X\beta^*, u_\ell \rangle}{\sigma} \right) \right), \end{aligned}$$

for $\ell \in \{D_{m^*} + 1, \dots, r\}$:

$$\begin{aligned} G_\ell &= 2 \left(1 - \Phi(\sqrt{\ell K}) \right) \\ H_\ell &= 2 \left(\Phi(\sqrt{\ell K}) - \Phi(\sqrt{K}) \right), \end{aligned}$$

3. $\forall r \in \{D_{m^*} + 1, \dots, q\}$, $\bar{f}_r(\cdot, \beta^*, \sigma^2)$ is defined by:

$$\begin{aligned} \bar{f}_r(K, \beta^*, \sigma^2) &= 1 - \max \left(\max_{\ell \in \{1, \dots, r-D_{m^*}\}} \left(F_{\chi^2(\ell)}(\ell K) \right), \right. \\ &\quad \left. \max_{\ell \in \{r-D_{m^*}+1, \dots, r\}} \left(F_{\chi^2(\ell)} \left(\frac{\ell K}{2} - \sum_{k=r-\ell+1}^{D_{m^*}} \frac{\langle X\beta^*, u_k \rangle^2}{\sigma^2} \right) \right) \right). \end{aligned}$$

Proof of Theorem 3.2 can be found in Subsection 6.2.

Hence, although the model selection procedure is built for prediction performances, bounds on the FDR are derived with respect to $\hat{m}(K)$. Terms $\underline{f}_r(K, \beta^*, \sigma^2)$ and $\bar{f}_r(K, \beta^*, \sigma^2)$ only involve some evaluations of cumulative distribution functions of the standard Gaussian and chi-squared variables. So, they have a fully explicit form which makes easier the understanding of the behavior of the FDR in model selection. Note that some of these terms depend on the signal-to-noise ratio, as usual in statistics.

3.2.3 Strictly positive FDR.

The following corollary gives a lower bound on the FDR which is independent from σ^2 .

Corollary 3.3. *Under the assumptions and definitions of Theorem 3.2, $\forall K > 0$:*

$$FDR(\hat{m}(K)) \geq \sum_{r=D_{m^*}+1}^q \left(\frac{r - D_{m^*}}{r} P_r(K) \frac{2\sqrt{2}}{\sqrt{\pi}(\sqrt{rK} + \sqrt{rK+4})} e^{-\frac{rK}{2}} \right) > 0.$$

Proof of Corollary 3.3 can be found in Subsection 6.3.

Hence, $FDR(\hat{m}(K)) > 0$ for all $K > 0$ and whatever σ^2 . This is not surprising since the $FDR(\hat{m}(K))$ is strictly positive even in the simplest case of no noise level. Indeed, when $\sigma^2 = 0$, $Y = X\beta^*$ and the minimization in (2.2) is reduced to the least-squares minimization. So, $\hat{\beta}_{m^*} = \beta^*$, which provides a zero value of the associated least-squares criterion. Moreover, for $m \in \mathcal{M}$ such that $m^* \subset m$, $\hat{\beta}_m$ leads to a zero value of the least-squares criterion with a non-zero probability. So, the selected model \hat{m} is strictly larger than m^* with a non-zero probability leading to a strictly positive value of the $FDR(\hat{m})$. By taking the expectation, $FDR(\hat{m}) > 0$. The larger σ^2 , the larger the FDR. So for $\sigma^2 > 0$, the event $m^* \subset \hat{m}$ happens with a non-zero probability providing a strictly positive FDR as well.

3.2.4 Asymptotic analysis.

The following corollary gives the asymptotic behavior of the FDR function in model selection when K tends to infinity.

Corollary 3.4. *Under the assumptions and the definitions of Theorem 3.2, the $FDR(\hat{m}(K))$ function tends to 0 when K tends to infinity and satisfies $\forall \eta > 0$,*

$$FDR(\hat{m}(K)) = o_{K \rightarrow +\infty} \left(e^{-K(\frac{1}{2}-\eta)} \right). \quad (3.6)$$

Furthermore, $\forall \eta > 0$, $\exists C_\eta > 0$, $\exists L_\eta > 0$, $\forall K > L_\eta$, we have:

$$FDR(\hat{m}(K)) \geq C_\eta e^{-K \left(\frac{D_{m^*}+1+2\eta}{2} \right)}. \quad (3.7)$$

So, $\forall \varepsilon > 0$,

$$\begin{aligned} -\frac{D_{m^*}}{2} - \frac{1}{2} - \varepsilon &\leq \liminf_{K \rightarrow +\infty} \frac{1}{K} \log \left(FDR(\hat{m}(K)) \right) \\ \limsup_{K \rightarrow +\infty} \frac{1}{K} \log \left(FDR(\hat{m}(K)) \right) &\leq -\frac{1}{2} + \varepsilon. \end{aligned} \quad (3.8)$$

Proof of Corollary 3.4 can be found in Subsection 6.4.

From Equation (3.6), $FDR(\hat{m}(K))$ tends to 0 when K tends to $+\infty$ with at least an exponential convergence rate and Equation (3.7) suggests that the exponential convergence rate is optimal. Moreover, although this result is asymptotic in K , Equation (3.6) states that there is no need to go far from $K = 2$ to have a reasonably small control of the FDR in the model selection procedure (1.4).

Remark 3.5. *With no signal ($\beta^* = 0$ and $D_{m^*} = 0$), the asymptotic bounds in (3.8) are $-\frac{1}{2} - \varepsilon$ and $-\frac{1}{2} + \varepsilon$ and consequently:*

$$\log \left(FDR(\hat{m}(K)) \right) \underset{K \rightarrow +\infty}{\sim} -\frac{1}{2}K.$$

Remark 3.6. *The asymptotic upper and lower bounds (3.6) and (3.7) are satisfied whatever the value of $\sigma^2 > 0$. It is possible to obtain the following sharpest asymptotic upper bound: $\forall \tilde{\eta} > 0$,*

$$FDR(\hat{m}(K)) = o \left(e^{-\left(K \frac{(D_{m^*}+1-\tilde{\eta})}{4} - \frac{1}{2\sigma^2} \sum_{k=1}^{D_{m^*}} \langle X\beta^*, u_k \rangle^2 \right)} \right) \quad (3.9)$$

in the asymptotic regime where $K \rightarrow +\infty$ and $\sigma \rightarrow 0$ with $\frac{1}{\sigma} = o_{\sigma \rightarrow 0}(\sqrt{K})$. The reader can find the proof in Subsection 6.4.

3.3 Illustrations of the main result in the orthogonal case.

We propose to analyze the particular case where the design matrix X is orthogonal since it leads to simplified forms for the FDR bounds easy to implement.

Corollary 3.7 (Application on the orthogonal case). *Under assumptions of Theorem 3.2 and by assuming that (X_1, \dots, X_q) are orthonormal for $\langle \cdot, \cdot \rangle$, then: $\forall K > 0$, $FDR(\hat{m}(K))$ satisfies the same inequalities as (3.4) where: for $\ell \in \{1, \dots, D_{m^*}\}$:*

$$G_\ell = 2 - \left(\Phi\left(\sqrt{\ell K} - \frac{\beta_\ell^*}{\sigma}\right) + \Phi\left(\sqrt{\ell K} + \frac{\beta_\ell^*}{\sigma}\right) \right),$$

for $\ell \in \{2, \dots, D_{m^*}\}$:

$$H_\ell = \Phi\left(\sqrt{\ell K} - \frac{\beta_\ell^*}{\sigma}\right) + \Phi\left(\sqrt{\ell K} + \frac{\beta_\ell^*}{\sigma}\right) - \left(\Phi\left(\sqrt{K} - \frac{\beta_\ell^*}{\sigma}\right) + \Phi\left(\sqrt{K} + \frac{\beta_\ell^*}{\sigma}\right) \right),$$

for all $r \in \{D_{m^*} + 1, \dots, q\}$:

$$\bar{f}_r(K, \beta^*, \sigma^2) = 1 - \max \left(\max_{\ell \in \{1, \dots, r - D_{m^*}\}} \left(F_{\chi^2(\ell)}(\ell K) \right), \max_{\ell \in \{r - D_{m^*} + 1, \dots, r\}} \left(F_{\chi^2(\ell)}\left(\frac{\ell K}{2} - \sum_{k=r-\ell+1}^{D_{m^*}} \frac{\beta_k^{*2}}{\sigma^2} \right) \right) \right),$$

and all other terms are the same as those defined in Theorem 3.2.

Proof of Corollary 3.7 can be found in Subsection 6.5.

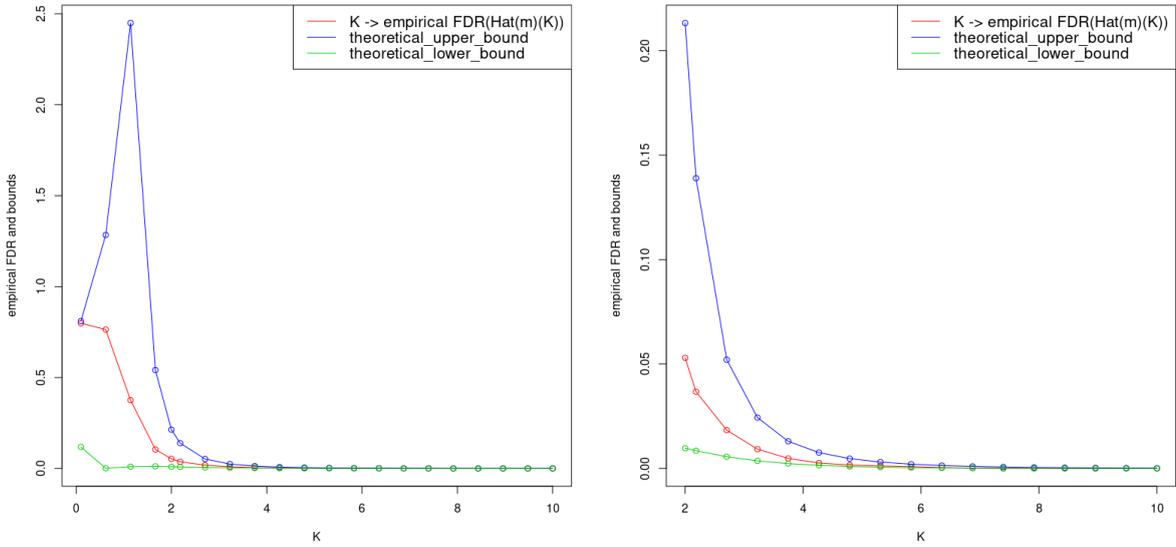


Figure 2: Left: curves of the empirical estimation of the $FDR(\hat{m}(K))$ (red) and the terms $b(K, \beta^*, \sigma^2)$ (green) and $B(K, \beta^*, \sigma^2)$ (blue) under the orthogonal design matrix X for the toy data set described in Subsection 7.1. Right: curves are plotted only for $K \geq 2$.

In Figure 2, we plot the empirical estimation of the $FDR(\hat{m}(K))$ with the functions $b(K, \beta^*, \sigma^2)$ and $B(K, \beta^*, \sigma^2)$ on a grid of positive K (left) and for $K \geq 2$ (right). Graphs are obtained from the *toy data set* described in Section 7 where X is orthogonal. The left figure is devoted to illustrate Corollary 3.7. The FDR function curve is well within the lower and upper bounds curves. From the right figure and consistency to Corollary 3.4, the empirical $FDR(\hat{m}(K))$ approaches 0 when K increases and the convergence rate seems to be exponential. Moreover, the curves of $b(K, \beta^*, \sigma^2)$ and $B(K, \beta^*, \sigma^2)$ frame the empirical FDR and the difference between the three curves becomes quickly negligible for K larger than 2.

4 Trade-off between the PR and the FDR controls.

While bounds $b(K, \beta^*, \sigma^2)$ and $B(K, \beta^*, \sigma^2)$ are easily understandable and fully implementable, they depend on β^* and σ^2 , unknown in practice. For a practical use, we propose to replace the theoretical bounds on the FDR as well as the theoretical expression of the PR with observable quantities from the data set (Subsection 4.1). Then, we propose an algorithm to calibrate the hyperparameter K from the data set such that both PR and FDR are small (Subsection 4.2).

4.1 Estimation of the theoretical terms.

When n is fixed, empirical approaches are not adapted. An alternative is to replace the theoretical terms by observable quantities.

4.1.1 Estimation of the PR.

Commonly, the predictive risk is evaluated with the mean squared error on a validation set independent from the training set used to estimate the parameters (see Formula (7.1) for the definition). However, it requires separating the dataset in two parts which increases the estimation errors. Here, we propose to use the entire dataset to both apply the model selection procedure and evaluate the predictive risk. By re-express the PR, we obtain the following proposition:

Proposition 4.1. *The dynamics of $\mathbb{E}[\|X\hat{\beta}_{\hat{m}(2)} - X\hat{\beta}_{\hat{m}(K)}\|_2^2]$ with respect to $K > 0$ is close to the one of $\mathbb{E}[\|Y - X\hat{\beta}_{\hat{m}(K)}\|_2^2]$.*

Proof of Proposition 4.1 can be found in Subsection 6.6.

Hence, the constant K minimizing $\mathbb{E}[\|X\hat{\beta}_{\hat{m}(2)} - X\hat{\beta}_{\hat{m}(K)}\|_2^2]$ and the one minimizing $\mathbb{E}[\|Y - X\hat{\beta}_{\hat{m}(K)}\|_2^2]$ are almost equal. Therefore, to evaluate the prediction performances, we propose the following term that we call *estimated risk*:

$$\widehat{\text{PR}}(\hat{m}(K)) = \frac{1}{n} \sum_{i=1}^n \left((X\hat{\beta}_{\hat{m}(2)})_i - (X\hat{\beta}_{\hat{m}(K)})_i \right)^2. \quad (4.1)$$

4.1.2 Estimation of the FDR.

The functions $b(\cdot, \beta^*, \sigma^2)$ and $B(\cdot, \beta^*, \sigma^2)$ are explicit and easily implementable but depend on β^* and σ^2 , both unknown.

We propose:

1. to apply the slope heuristic method [12] to get an estimator $\hat{\sigma}^2$ of σ^2 ,
2. to replace β^* by the estimator $\hat{\beta}_{\hat{m}(4)}$.

Justifications for the choice of these two estimates are provided in Section 7.

4.2 A data-dependent calibration of K in model selection procedure.

We propose a completely data-driven calibration of the hyperparameter K using the estimated risk function given in (4.1) and the $B(\cdot, \hat{\beta}_{\hat{m}(4)}, \hat{\sigma}^2)$ function to obtain a low value of both PR and FDR.

We propose the following algorithm:

Algorithm 1: Algorithm to calibrate K

1. Choose α the threshold for the FDR control.
 2. Compute $I_1 = \{K \geq 2, B(K, \hat{\beta}_{\hat{m}(4)}, \hat{\sigma}^2) \in]0, \alpha[\}$.
 3. Compute $I_2 = \{K \geq 2, \widehat{\text{PR}}(\hat{m}(K)) \approx \widehat{\text{PR}}(\hat{m}(2)) \}$.
 4. If $I_1 \cap I_2 \neq \emptyset$, return $\min \{K, K \in I_1 \cap I_2 \}$;
Otherwise, return $\min \{K, K \in I_1 \}$ or take a larger value of α .
-

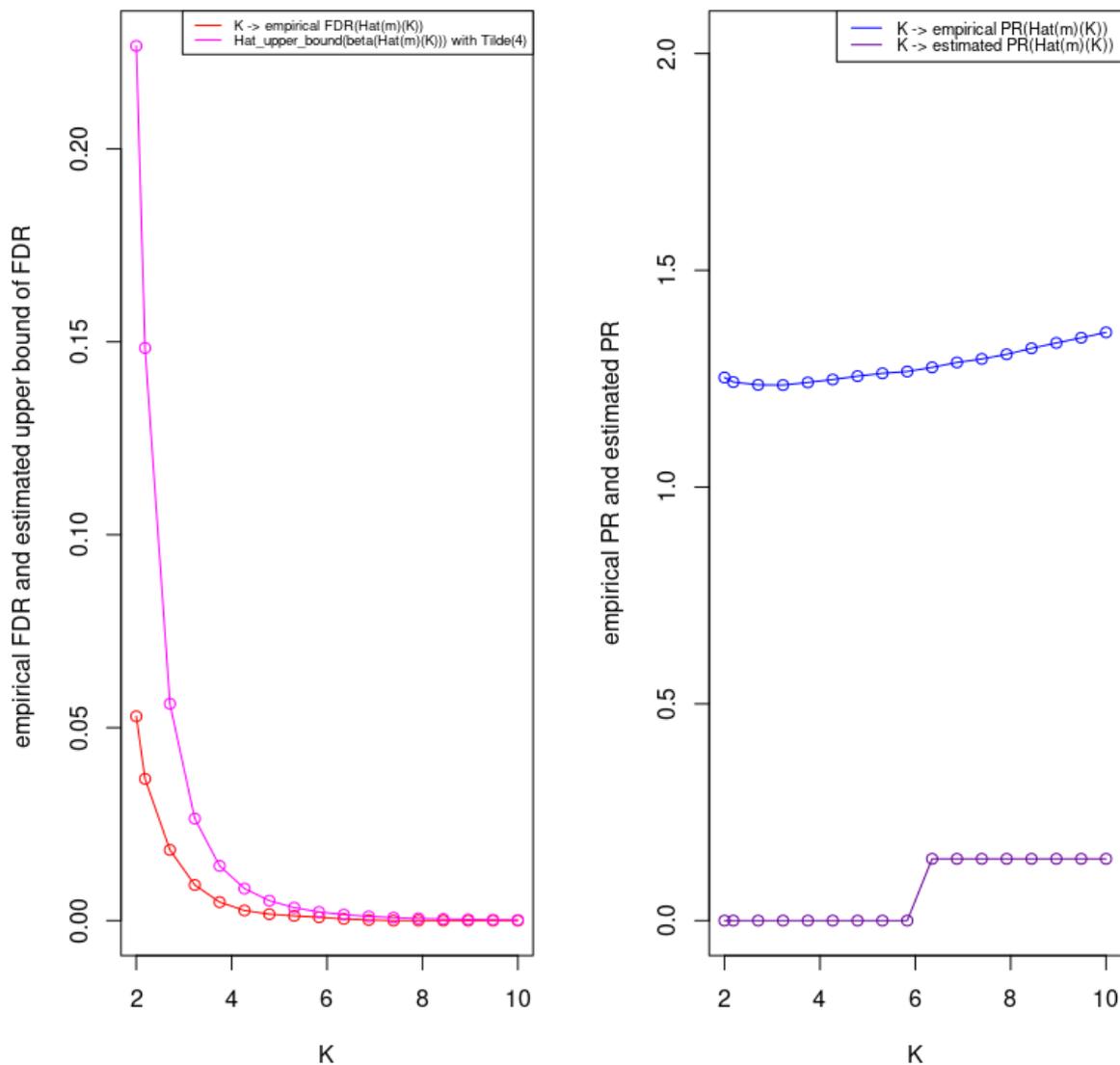


Figure 3: Curves of the empirical estimation functions $\widehat{\text{FDR}}(\hat{m}(K))$ (red) and $\widehat{\text{PR}}(\hat{m}(K))$ (blue), of the estimated risk (violet) and of the $B(K, \hat{\beta}_{\hat{m}(4)}, \hat{\sigma}^2)$ function (pink) for $K \geq 2$ for the toy data set described in Subsection 7.1.

Curves of Figure 3 are generated from the *toy data set* described in Section 7. For this example, we choose $\alpha = 0.05$ and the condition $|\widehat{\text{PR}}(\hat{m}(K)) - \widehat{\text{PR}}(\hat{m}(2))| \leq 0.1$ for the definition of $\widehat{\text{PR}}(\hat{m}(K)) \approx \widehat{\text{PR}}(\hat{m}(2))$ in the algorithm. We get $I_1 = [3.3, 10]$ and $I_2 = [2, 5.8]$ and so, our proposed algorithm returns $K = 3.3$. The selected model $\hat{m}(3.3)$ satisfies $\widehat{\text{PR}}(\hat{m}(3.3)) = 1.14$ and $B(3.3, \hat{\beta}_{\hat{m}(4)}, \hat{\sigma}^2) = 0.03$. This constant corresponds to a low value of both empirical predictive risk and FDR curves. Indeed, the empirical predictive risk of $\hat{m}(3.3)$ is equal to 1.24 and the empirical FDR of $\hat{m}(3.3)$ is equal to 0.01. To compare with the usual choice $K = 2$, the empirical predictive risk of $\hat{m}(3.3)$ is equal to 1.25 and the empirical FDR of $\hat{m}(3.3)$ is equal to 0.05. Hence, our proposed algorithm allows to maintain the prediction performances from $\hat{m}(2)$, reinforce the control of the FDR criterion and so gain a convenient trade-off between PR and FDR.

In Section 7, this algorithm 1 is applied on several data sets generated from various sets of parameters. Each time, the hyperparameter K is strictly larger than the commonly used constant 2.

5 Conclusions.

The variable selection procedure in a high-dimensional Gaussian linear regression with sparsity assumption is commonly used to identify a set of variables with prediction performances or to avoid the selection of non active variables. For prediction performances, the PR is usually controlled via a penalized least-squares minimization; to avoid the selection of non active variables, the FDR is usually controlled via a multiple testing approach. Controlling the PR tends to select too many variables, including non active ones, whereas controlling the FDR tends to select too few variables, leaving out some active ones.

This work shows that a convenient trade-off between PR and FDR can be achieved in ordered variable selection. The originality of this paper is to obtain this trade-off through a proper calibration of the hyperparameter K in the penalty of the model selection (1.4). Firstly, theoretical results lead to non-asymptotic lower and upper bounds on the FDR($\hat{m}(K)$) function when σ^2 is known. Asymptotic behaviors suggest that bounds are optimal. Secondly, the proposed methodology provides an algorithm to calibrate the hyperparameter K in the penalty function when σ^2 is unknown. This algorithm is based on completely data-driven terms: the estimated risk and the estimated upper bound on the FDR where the choices of estimators $\hat{\sigma}^2$ and $\hat{\beta}_{\hat{m}(4)}$ are derived from an extensive simulation study. The hyperparameter K is calibrated from the dataset to ensure $\widehat{\text{PR}}(\hat{m}(K)) \approx \widehat{\text{PR}}(\hat{m}(2))$ under the constraint $B(K, \hat{\beta}_{\hat{m}(4)}, \hat{\sigma}^2) < \alpha$. Our algorithm is validated on an extensive simulation study and allows to obtain a selected model ensuring a small value of both theoretical PR and FDR. The calibrated hyperparameter K is strictly larger than the commonly used constant $K = 2$.

If $D_{\hat{m}(K)} = q$ for one $K > 1$, the lower and upper bounds equal 0. This means that if $D_{\hat{m}(K)} = q$, a distinction between $D_{m^*} = q$ and $D_{m^*} < q$ is not possible without additional arguments. This is a limitation of our work.

To establish Theorem 3.2, variables are supposed to be ordered. The main perspective of our work is to generalize this result to complete variable selection procedure. This requires a complex combinatorial computation that appears at the stage of formula (3.2) as well as the use of a more complicated form of penalty function in model selection. Another generalization is to random model collections or a non-fixed design, more general frameworks adapted to some application points of view. These extensions are much more intricate. A second perspective is the construction of other algorithms based on the theoretical results. Indeed, estimators of σ^2 and β^* to get observable quantities and our proposed calibration of K are not necessarily the optimal choices. For example, an idea is to rerun Algorithm 1 several times, updating the value of K for the β^* estimate of the FDR bound by the algorithm output. The algorithm would then be less sensitive to the choice of the proposed $\hat{\beta}_{\hat{m}(4)}$ as input. A possible opening is to study the potential characteristics of the hyperparameter K provided by our data-driven hyperparameter calibration in a theoretical point of view, for instance K could depend on $|\beta^*|$. Finally, another possible extension is to study the false negative rate (FNR) function in the model selection procedure, similarly and in addition to the FDR one. This can be increase the power of the method, in the same idea of [31, 32].

6 Proofs of theoretical results.

This section contains proofs of all the theoretical results of this paper.

6.1 FDR expression in model selection.

Proof of Formula 3.1.

If $D_m^* = q$, then $\text{FP}(m) = 0$ for all $m \in \mathcal{M}$ and $\text{FDR}(m) = 0$ for all $m \in \mathcal{M}$.

Let us now suppose that $D_m^* < q$. The FDP expression within the model selection procedure is:

$$\begin{aligned}
 \forall K > 0, \quad \text{FDP}(\hat{m}(K)) &= \frac{\text{FP}(\hat{m}(K))}{\max(D_{\hat{m}(K)}, 1)} \\
 &\stackrel{(*)}{=} \frac{D_{\hat{m}(K)} - D_{m^*}}{D_{\hat{m}(K)}} \mathbb{1}_{\{D_{\hat{m}(K)} > D_{m^*}\}} \\
 &= \sum_{r=1}^q \frac{r - D_{m^*}}{r} \mathbb{1}_{\{r > D_{m^*}\}} \mathbb{1}_{\{D_{\hat{m}(K)} = r\}} \\
 &\stackrel{(**)}{=} \sum_{r=D_{m^*}+1}^q \frac{r - D_{m^*}}{r} \mathbb{1}_{\{\hat{m}(K) = m_r\}} \\
 &\stackrel{(***)}{=} \sum_{r=D_{m^*}+1}^q \frac{r - D_{m^*}}{r} \mathbb{1}_{\left\{ \bigcap_{\substack{\ell=0 \\ \ell \neq r}}^q \{\text{crit}_K(m_r) < \text{crit}_K(m_\ell)\} \right\}}.
 \end{aligned}$$

(*) and (**) are due to the fact that models $(m)_{m \in \mathcal{M}}$ are nested and $m^* \in \mathcal{M}$. (***) is obtained since the crit_K function is injective on \mathcal{M} . Finally, by taking the expectation, we obtain the FDR expression (3.1). \square

Proof of Proposition 3.1.

Before proving Proposition 3.1, let us cite and prove two lemmas.

Lemma 6.1. For $r \in \{D_{m^*} + 1, \dots, q\}$ and for all $\ell \in \{0, \dots, r - 1\}$:

$$\|Y - X\hat{\beta}_{m_r}\|_2^2 - \|Y - X\hat{\beta}_{m_\ell}\|_2^2 = - \sum_{k=\ell+1}^r \langle Y, u_k \rangle^2.$$

Lemma 6.2. For $r \in \{D_{m^*} + 1, \dots, q\}$ and for all $\ell \in \{r + 1, \dots, q\}$:

$$\|Y - X\hat{\beta}_{m_r}\|_2^2 - \|Y - X\hat{\beta}_{m_\ell}\|_2^2 = \sum_{k=r+1}^{\ell} \langle Y, u_k \rangle^2.$$

Proof of Lemma 6.1.

For $r \in \{D_{m^*} + 1, \dots, q\}$ and $\ell \in \{0, \dots, r - 1\}$:

$$\begin{aligned}
 \|Y - X\hat{\beta}_{m_r}\|_2^2 - \|Y - X\hat{\beta}_{m_\ell}\|_2^2 &= \|X\hat{\beta}_{m_r}\|_2^2 - \|X\hat{\beta}_{m_\ell}\|_2^2 + 2\langle Y, X\hat{\beta}_{m_\ell} - X\hat{\beta}_{m_r} \rangle \\
 &= \|X\hat{\beta}_{m_r}\|_2^2 - \|X\hat{\beta}_{m_\ell}\|_2^2 + 2\langle Y - X\hat{\beta}_{m_r}, X\hat{\beta}_{m_\ell} \rangle \\
 &\quad - 2\langle Y - X\hat{\beta}_{m_r}, X\hat{\beta}_{m_r} \rangle + 2\langle X\hat{\beta}_{m_r}, X\hat{\beta}_{m_\ell} \rangle - 2\|X\hat{\beta}_{m_r}\|_2^2 \\
 &= -\|X\hat{\beta}_{m_r}\|_2^2 - \|X\hat{\beta}_{m_\ell}\|_2^2 + 2\langle X\hat{\beta}_{m_r}, X\hat{\beta}_{m_\ell} \rangle = -\|X\hat{\beta}_{m_r} - X\hat{\beta}_{m_\ell}\|_2^2.
 \end{aligned}$$

The last line is due to the fact that $Y - X\hat{\beta}_{m_r} \in (m_r)^\perp \subset (m_\ell)^\perp$ since $m_\ell \subset m_r$ and $X\hat{\beta}_{m_r}$ is the projection of Y onto m_r .

Then,

$$\begin{aligned}
 \|X\hat{\beta}_{m_r} - X\hat{\beta}_{m_\ell}\|_2^2 &= \|\Pi_{m_r}(Y) - \Pi_{m_\ell}(Y)\|_2^2 \\
 &= \|\Pi_{\text{Span}(X_1, \dots, X_r)}(Y) - \Pi_{\text{Span}(X_1, \dots, X_\ell)}(Y)\|_2^2 \\
 &\stackrel{(*)}{=} \|\Pi_{\text{Span}(u_1, \dots, u_r)}(Y) - \Pi_{\text{Span}(u_1, \dots, u_\ell)}(Y)\|_2^2 \\
 &= \|\Pi_{\text{Span}(u_{\ell+1}, \dots, u_r)}(Y)\|_2^2 \\
 &= \left\| \sum_{k=\ell+1}^r \langle Y, u_k \rangle u_k \right\|_2^2 \\
 &\stackrel{(**)}{=} \sum_{k=\ell+1}^r \langle Y, u_k \rangle^2.
 \end{aligned}$$

(*) come from the definition of (u_1, \dots, u_n) and (**) is obtained by Parseval's identity. \square

Proof of Lemma 6.2.

For $r \in \{D_{m^*} + 1, \dots, q\}$ and $\ell \in \{r + 1, \dots, q\}$:

$$\begin{aligned}
 \|Y - X\hat{\beta}_{m_r}\|_2^2 - \|Y - X\hat{\beta}_{m_\ell}\|_2^2 &= \|X\hat{\beta}_{m_r}\|_2^2 - \|X\hat{\beta}_{m_\ell}\|_2^2 + 2\langle Y, X\hat{\beta}_{m_\ell} - X\hat{\beta}_{m_r} \rangle \\
 &= \|X\hat{\beta}_{m_r}\|_2^2 - \|X\hat{\beta}_{m_\ell}\|_2^2 + 2\langle Y - X\hat{\beta}_{m_\ell}, X\hat{\beta}_{m_\ell} \rangle \\
 &\quad - 2\langle Y - X\hat{\beta}_{m_\ell}, X\hat{\beta}_{m_r} \rangle + 2\|X\hat{\beta}_{m_\ell}\|_2^2 \\
 &\quad - 2\langle X\hat{\beta}_{m_\ell}, X\hat{\beta}_{m_r} \rangle \\
 &\stackrel{(*)}{=} \|X\hat{\beta}_{m_r}\|_2^2 + \|X\hat{\beta}_{m_\ell}\|_2^2 - 2\langle X\hat{\beta}_{m_\ell}, X\hat{\beta}_{m_r} \rangle \\
 &= \|X\hat{\beta}_{m_\ell} - X\hat{\beta}_{m_r}\|_2^2.
 \end{aligned}$$

(*) is due to the fact that $Y - X\hat{\beta}_{m_\ell} \in (m_\ell)^\perp \subset (m_r)^\perp$ since $m_r \subset m_\ell$, and $X\hat{\beta}_{m_\ell}$ is the projection of Y onto m_ℓ . Then,

$$\begin{aligned}
 \|X\hat{\beta}_{m_\ell} - X\hat{\beta}_{m_r}\|_2^2 &= \|\Pi_{m_\ell}(Y) - \Pi_{m_r}(Y)\|_2^2 \\
 &= \|\Pi_{\text{Span}(X_1, \dots, X_\ell)}(Y) - \Pi_{\text{Span}(X_1, \dots, X_r)}(Y)\|_2^2 \\
 &\stackrel{(*)}{=} \|\Pi_{\text{Span}(u_1, \dots, u_\ell)}(Y) - \Pi_{\text{Span}(u_1, \dots, u_r)}(Y)\|_2^2 \\
 &= \|\Pi_{\text{Span}(u_{r+1}, \dots, u_\ell)}(Y)\|_2^2 \\
 &= \left\| \sum_{k=r+1}^{\ell} \langle Y, u_k \rangle u_k \right\|_2^2 \\
 &\stackrel{(**)}{=} \sum_{k=r+1}^{\ell} \langle Y, u_k \rangle^2.
 \end{aligned}$$

(*) come from the definition of (u_1, \dots, u_n) and (**) is obtained by Parseval's identity. \square

Proof of Proposition 3.1.

Starting from (3.1), we decompose the event $\left\{ \bigcap_{\substack{\ell=0 \\ \ell \neq r}}^q \{\text{crit}_K(m_r) < \text{crit}_K(m_\ell)\} \right\}$ by the intersection of these two events

$$\left\{ \bigcap_{\ell=0}^{r-1} \{\text{crit}_K(m_r) < \text{crit}_K(m_\ell)\} \right\} \text{ and } \left\{ \bigcap_{\ell=r+1}^q \{\text{crit}_K(m_r) < \text{crit}_K(m_\ell)\} \right\}.$$

By using the definition of the crit_K function, we have for $r \in \{D_{m^*} + 1, \dots, q\}$ and $\ell \in \{0, \dots, r-1, r+1, \dots, q\}$:

$$\begin{aligned}
 \left\{ \text{crit}_K(m_r) < \text{crit}_K(m_\ell) \right\} &= \left\{ \|Y - X\hat{\beta}_{m_r}\|_2^2 + K\sigma^2 r < \|Y - X\hat{\beta}_{m_\ell}\|_2^2 + K\sigma^2 \ell \right\} \\
 &= \left\{ \|Y - X\hat{\beta}_{m_r}\|_2^2 - \|Y - X\hat{\beta}_{m_\ell}\|_2^2 < K\sigma^2(\ell - r) \right\}.
 \end{aligned}$$

So, by applying Lemma 6.1, $\ell \in \{0, \dots, r-1\}$:

$$\left\{ \text{crit}_K(m_r) < \text{crit}_K(m_\ell) \right\} = \left\{ \sum_{k=\ell+1}^r \langle Y, u_k \rangle^2 > K\sigma^2(r-\ell) \right\},$$

and by applying Lemma 6.2, $\ell \in \{r+1, \dots, q\}$:

$$\left\{ \text{crit}_K(m_r) < \text{crit}_K(m_\ell) \right\} = \left\{ \sum_{k=r+1}^{\ell} \langle Y, u_k \rangle^2 < K\sigma^2(\ell-r) \right\}.$$

In this way, $\left\{ \bigcap_{\substack{\ell=0 \\ \ell \neq r}}^q \left\{ \text{crit}_K(m_r) < \text{crit}_K(m_\ell) \right\} \right\}$ is decomposed by two events:

$$\left\{ \bigcap_{\ell=0}^{r-1} \left\{ \sum_{k=\ell+1}^r \langle Y, u_k \rangle^2 > K\sigma^2(r-\ell) \right\} \right\} \cap \left\{ \bigcap_{\ell=r+1}^q \left\{ \sum_{k=r+1}^{\ell} \langle Y, u_k \rangle^2 < K\sigma^2(\ell-r) \right\} \right\}.$$

Let us define U the $n \times n$ matrix such that u_k is the k -th column of U . Since $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ and (u_1, \dots, u_n) is an orthonormal basis of \mathbb{R}^n , we get $U^T \varepsilon = (\langle \varepsilon, u_1 \rangle, \dots, \langle \varepsilon, u_n \rangle)^T \sim \mathcal{N}(0, \sigma^2 U I_n U^T) = \mathcal{N}(0, \sigma^2 I_n)$. Hence, random variables $(\langle Y, u_i \rangle)_{i \in \{1, \dots, n\}}$ are independent with $\langle Y, u_i \rangle \sim \mathcal{N}(\langle X\beta^*, u_i \rangle, \sigma^2)$ for all i in $\{1, \dots, n\}$. Since the first event of the previous decomposition depends only on random variables $\langle Y, u_i \rangle$ for $i \in \{1, \dots, r-1\}$ whereas the second one depends only on random variables $\langle Y, u_i \rangle$ for $i \in \{r+1, \dots, q\}$, the two events are independent. Hence, from (3.1), we obtain for all $K > 0$:

$$\begin{aligned} \text{FDR}(\hat{m}(K)) &= \sum_{r=D_{m^*}+1}^q \frac{r-D_{m^*}}{r} \mathbb{P} \left(\bigcap_{\ell=0}^{r-1} \left\{ \sum_{k=\ell+1}^r \langle Y, u_k \rangle^2 > K\sigma^2(r-\ell) \right\} \right) \\ &\quad \times \mathbb{P} \left(\bigcap_{\ell=r+1}^q \left\{ \sum_{k=r+1}^{\ell} \langle Y, u_k \rangle^2 < K\sigma^2(\ell-r) \right\} \right). \end{aligned}$$

Moreover, since $\langle X\beta^*, u_k \rangle = 0, \forall k > D_{m^*}$ and since $r \geq D_{m^*} + 1$, we have:

$$\sum_{k=\ell+1}^r \langle Y, u_k \rangle^2 = \sum_{k=\ell+1}^r \langle \varepsilon, u_k \rangle^2.$$

So, for all $K > 0$ and for each $r \in \{D_{m^*} + 1, \dots, q\}$:

$$\begin{aligned} \mathbb{P} \left(\bigcap_{\ell=r+1}^q \left\{ \sum_{k=r+1}^{\ell} \langle Y, u_k \rangle^2 < K\sigma^2(\ell-r) \right\} \right) &= \mathbb{P} \left(\bigcap_{\ell=r+1}^q \left\{ \sum_{k=r+1}^{\ell} \tilde{Z}_k^2 < K\sigma^2(\ell-r) \right\} \right), \\ &\quad \text{where } \tilde{Z}_k \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2) \\ \mathbb{P} \left(\bigcap_{\ell=r+1}^q \left\{ \sum_{k=r+1}^{\ell} \langle Y, u_k \rangle^2 < K\sigma^2(\ell-r) \right\} \right) &= \mathbb{P} \left(\bigcap_{\ell=r+1}^q \left\{ \sum_{k=r+1}^{\ell} Z_k^2 < K(\ell-r) \right\} \right), \\ &\quad \text{where } Z_k \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1). \end{aligned}$$

Hence, for all $K > 0$ and for each $r \in \{D_{m^*} + 1, \dots, q\}$,

$\mathbb{P} \left(\bigcap_{\ell=r+1}^q \left\{ \sum_{k=r+1}^{\ell} \langle Y, u_k \rangle^2 < K\sigma^2(\ell-r) \right\} \right)$ does not depend on the data and we deduce the Formula (3.2) with:

$$\begin{aligned} P_r(K) &= \mathbb{P} \left(\bigcap_{\ell=r+1}^q \left\{ \sum_{k=r+1}^{\ell} Z_k^2 < K(\ell-r) \right\} \right), \\ Q_r(K, \beta^*, \sigma^2) &= \mathbb{P} \left(\bigcap_{\ell=0}^{r-1} \left\{ \sum_{k=\ell+1}^r \langle Y, u_k \rangle^2 > K\sigma^2(r-\ell) \right\} \right), \end{aligned}$$

where $Z_k \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1), \forall k \in \{r+1, \dots, q\}$. □

6.2 General bounds.

Proof of Theorem 3.2.

We start from (3.2).

- bounds on the Q_r terms.

For all $K > 0$ and for each $r \in \{D_{m^*} + 1, \dots, q\}$, we recall that:

$$Q_r(K, \beta^*, \sigma^2) = \mathbb{P} \left(\bigcap_{\ell=0}^{r-1} \left\{ \sum_{k=\ell+1}^r \langle Y, u_k \rangle^2 > K\sigma^2(r - \ell) \right\} \right),$$

and since $\langle X\beta^*, u_k \rangle = 0, \forall k > D_{m^*}$, we have:

$$\begin{aligned} Q_r(K, \beta^*, \sigma^2) &= \mathbb{P} \left(\bigcap_{\ell=0}^{r-1} \left\{ \sum_{k=\ell+1}^r \left(\langle \varepsilon, u_k \rangle^2 \mathbf{1}_{k > D_{m^*}} + \langle Y, u_k \rangle^2 \mathbf{1}_{k \leq D_{m^*}} \right) > K\sigma^2(r - \ell) \right\} \right) \\ &= \mathbb{P} \left(\left\{ \langle \varepsilon, u_r \rangle^2 > K\sigma^2 \right\} \cap \dots \cap \left\{ \langle \varepsilon, u_r \rangle^2 + \dots + \langle \varepsilon, u_{D_{m^*}+1} \rangle^2 > K\sigma^2(r - D_{m^*}) \right\} \right. \\ &\quad \cap \left\{ \langle \varepsilon, u_r \rangle^2 + \dots + \langle \varepsilon, u_{D_{m^*}+1} \rangle^2 + \langle Y, u_{D_{m^*}} \rangle^2 > K\sigma^2(r - D_{m^*} + 1) \right\} \cap \dots \\ &\quad \left. \cap \left\{ \langle \varepsilon, u_r \rangle^2 + \dots + \langle \varepsilon, u_{D_{m^*}+1} \rangle^2 + \langle Y, u_{D_{m^*}} \rangle^2 + \dots + \langle Y, u_1 \rangle^2 > K\sigma^2 r \right\} \right) \\ &= \mathbb{P} \left(\left\{ c_r > K\sigma^2 \right\} \cap \left\{ c_r + c_{r-1} > 2K\sigma^2 \right\} \cap \dots \cap \left\{ c_r + c_{r-1} + \dots + c_1 > rK\sigma^2 \right\} \right) \quad (6.1) \end{aligned}$$

where $c_\ell = \langle Y, u_\ell \rangle^2$ for $\ell \in \{1, \dots, D_{m^*}\}$ and $c_\ell = \langle \varepsilon, u_\ell \rangle^2$ for $\ell \in \{D_{m^*} + 1, \dots, r\}$.

Lower bound on $Q_r(K, \beta^*, \sigma^2)$ for $r \in \{D_{m^*} + 1, \dots, q\}$:

Lemma 6.3. *Let us consider an integer $s > 1$, $K > 0$ and c_1, \dots, c_s s non-negative random independent quantities. We define by E_ℓ the event $\{c_\ell > \ell K\sigma^2\}$ for $\ell \in \{1, \dots, s\}$ and by F_ℓ the event $\{K\sigma^2 < c_\ell \leq \ell K\sigma^2\}$ for $\ell \in \{2, \dots, s\}$. Then:*

$$\begin{aligned} &\left\{ c_s > K\sigma^2 \right\} \cap \left\{ c_s + c_{s-1} > 2K\sigma^2 \right\} \cap \dots \cap \left\{ c_s + c_{s-1} + \dots + c_1 > sK\sigma^2 \right\} \\ &\supseteq E_s \sqcup \left(F_s \cap \left(E_{s-1} \sqcup \left(F_{s-1} \cap \left(E_{s-2} \sqcup \dots \sqcup \left(F_3 \cap \left(E_2 \sqcup \left(F_2 \cap E_1 \right) \right) \right) \right) \right) \right) \right), \end{aligned}$$

where \cap and \sqcap design respectively any intersection and a disjoint intersection of events, as well as \cup and \sqcup designing respectively any union and a disjoint union of events.

Proof. We prove Lemma 6.3 by a recurrence on $s \geq 1$.

For $s = 1$, both sets correspond to E_1 , so the inclusion is obvious. Let $s \geq 1$ and suppose that the inclusion is true for s . With the definitions of E_{s+1} and F_{s+1} , we obtain:

$$\begin{aligned}
& \left\{ c_{s+1} > K\sigma^2 \right\} \cap \left\{ c_{s+1} + c_s > 2K\sigma^2 \right\} \cap \cdots \cap \left\{ c_{s+1} + c_s + \cdots + c_1 > (s+1)K\sigma^2 \right\} \\
&= \left(E_{s+1} \sqcup F_{s+1} \right) \\
& \cap \left(\left\{ c_{s+1} + c_s > 2K\sigma^2 \right\} \cap \cdots \cap \left\{ c_{s+1} + c_s + \cdots + c_1 > (s+1)K\sigma^2 \right\} \right) \\
&= \left(E_{s+1} \cap \left(\left\{ c_{s+1} + c_s > 2K\sigma^2 \right\} \cap \cdots \cap \left\{ c_{s+1} + c_s + \cdots + c_1 > (s+1)K\sigma^2 \right\} \right) \right) \\
& \sqcup \left(F_{s+1} \cap \left(\left\{ c_{s+1} + c_s > 2K\sigma^2 \right\} \cap \cdots \cap \left\{ c_{s+1} + c_s + \cdots + c_1 > (s+1)K\sigma^2 \right\} \right) \right) \\
&\stackrel{(*)}{=} E_{s+1} \\
& \sqcup \left(F_{s+1} \cap \left(\left\{ c_{s+1} + c_s > 2K\sigma^2 \right\} \cap \cdots \cap \left\{ c_{s+1} + c_s + \cdots + c_1 > (s+1)K\sigma^2 \right\} \right) \right) \\
&\stackrel{(**)}{\supseteq} E_{s+1} \sqcup \left(F_{s+1} \right. \\
& \left. \cap \left(\left\{ c_s > K\sigma^2 \right\} \cap \left\{ c_s + c_{s-1} > 2K\sigma^2 \right\} \cap \cdots \cap \left\{ c_s + c_{s-1} + \cdots + c_1 > sK\sigma^2 \right\} \right) \right) \\
&\stackrel{(***)}{\supseteq} E_{s+1} \sqcup \left(F_{s+1} \cap \left(E_s \sqcup \left(F_s \cap \left(E_{s-1} \sqcup \cdots \sqcup \left(F_3 \cap \left(E_3 \sqcup \left(F_2 \cap E_1 \right) \right) \right) \right) \right) \right) \right) \\
&\stackrel{(*****)}{\supseteq} E_{s+1} \sqcup \left(F_{s+1} \cap \left(E_s \sqcup \left(F_s \cap \left(E_{s-1} \sqcup \cdots \sqcup \left(F_3 \cap \left(E_3 \sqcup \left(F_2 \cap E_1 \right) \right) \right) \right) \right) \right) \right).
\end{aligned}$$

(*) is true since c_i are non-negative for all $i \in \{1, \dots, s+1\}$ providing that

$E_{s+1} \subset \left(\left\{ c_{s+1} + c_s > 2K\sigma^2 \right\} \cap \cdots \cap \left\{ c_{s+1} + c_s + \cdots + c_1 > (s+1)K\sigma^2 \right\} \right)$, (**) comes from the inclusion $\left\{ c_{s+1} > K\sigma^2 \right\} \subset F_{s+1}$. We obtain (***) by applying the recurrence assumption at the step s . Independence of c_1, \dots, c_{s+1} provides the independence between F_{s+1} and $\left(E_s \sqcup \left(F_s \cap \left(E_{s-1} \sqcup \cdots \sqcup \left(F_3 \cap \left(E_3 \sqcup \left(F_2 \cap E_1 \right) \right) \right) \right) \right) \right)$ which gets (****).

Thus, the property is true for $s+1$, which proves lemma. \square

By applying Lemma 6.3 on Formula (6.1) with $s = r$, we obtain:

$$\begin{aligned}
& Q_r(K, \beta^*, \sigma^2) \geq \mathbb{P}(E_r) \\
& + \mathbb{P}(F_r) \left(\mathbb{P}(E_{r-1}) + \mathbb{P}(F_{r-1}) \left(\mathbb{P}(E_{r-2}) + \cdots + \mathbb{P}(F_3) \left(\mathbb{P}(E_2) + \mathbb{P}(F_2) \mathbb{P}(E_1) \right) \right) \right).
\end{aligned}$$

By using that $\langle Y, u_\ell \rangle \sim \mathcal{N}(\langle X\beta^*, u_\ell \rangle, \sigma^2)$ for $\ell \in \{1, \dots, D_{m^*}\}$ and $\langle \varepsilon, u_\ell \rangle \in \mathcal{N}(0, \sigma^2)$ for $\ell \in \{1, \dots, r\}$, we get: For $\ell \in \{1, \dots, D_{m^*}\}$:

$$\begin{aligned}
\mathbb{P}(E_\ell) &= \mathbb{P} \left(\left\{ \langle Y, u_\ell \rangle^2 > \ell K \sigma^2 \right\} \right) \\
&= 2 - \left(\Phi \left(\sqrt{\ell K} - \frac{\langle X\beta^*, u_\ell \rangle}{\sigma} \right) + \Phi \left(\sqrt{\ell K} + \frac{\langle X\beta^*, u_\ell \rangle}{\sigma} \right) \right) \\
&= G_\ell.
\end{aligned}$$

Trade-off between prediction and FDR for variable selection

For $\ell \in \{2, \dots, D_{m^*}\}$:

$$\begin{aligned} \mathbb{P}(F_\ell) &= \mathbb{P}\left(\left\{K\sigma^2 < \langle Y, u_\ell \rangle^2 \leq \ell K\sigma^2\right\}\right) \\ &= \Phi\left(\sqrt{\ell K} - \frac{\langle X\beta^*, u_\ell \rangle}{\sigma}\right) + \Phi\left(\sqrt{\ell K} + \frac{\langle X\beta^*, u_\ell \rangle}{\sigma}\right) \\ &\quad - \left(\Phi\left(\sqrt{K} - \frac{\langle X\beta^*, u_\ell \rangle}{\sigma}\right) + \Phi\left(\sqrt{K} + \frac{\langle X\beta^*, u_\ell \rangle}{\sigma}\right)\right) \\ &= H_\ell. \end{aligned}$$

For $\ell \in \{D_{m^*} + 1, \dots, r\}$:

$$\begin{aligned} \mathbb{P}(E_\ell) &= \mathbb{P}\left(\left\{\langle \varepsilon, u_\ell \rangle^2 > \ell K\sigma^2\right\}\right) \\ &= 2\left(1 - \Phi(\sqrt{\ell K})\right) \\ &= G_\ell, \\ \mathbb{P}(F_\ell) &= \mathbb{P}\left(\left\{K\sigma^2 < \langle \varepsilon, u_\ell \rangle^2 \leq \ell K\sigma^2\right\}\right) \\ &= 2\left(\Phi(\sqrt{\ell K}) - \Phi(\sqrt{K})\right) \\ &= H_\ell. \end{aligned}$$

Hence, a lower bound on $Q_r(K, \beta^*, \sigma^2)$ is obtained for all $K > 0$:

$$\underline{f}_r(K, \beta^*, \sigma^2) \leq Q_r(K, \beta^*, \sigma^2) \tag{6.2}$$

with:

$$\begin{aligned} \underline{f}_r(K, \beta^*, \sigma^2) &= G_r + H_r \underline{f}_{r-1}(K, \beta^*, \sigma^2) \\ \text{and } \underline{f}_1(K, \beta^*, \sigma^2) &= G_1. \end{aligned} \tag{6.3}$$

Upper bound on $Q_r(K, \beta^*, \sigma^2)$ for $r \in \{D_{m^*} + 1, \dots, q\}$:

By using definitions of Lemma 6.3 and formula (6.1), we get:

$$\begin{aligned} Q_r(K, \beta^*, \sigma^2) &\leq \min\left(\mathbb{P}\left(\left\{c_r > K\sigma^2\right\}\right), \mathbb{P}\left(\left\{c_r + c_{r-1} > 2K\sigma^2\right\}\right), \dots, \right. \\ &\quad \left. \mathbb{P}\left(\left\{c_r + c_{r-1} + \dots + c_1 > rK\sigma^2\right\}\right)\right). \end{aligned} \tag{6.4}$$

Since $\langle \varepsilon, u_i \rangle_{i \in \{D_{m^*} + 1, \dots, r\}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2)$, we have for all $j \in \{D_{m^*} + 1, \dots, r\}$:

$$\mathbb{P}\left(\left\{c_r + \dots + c_j > (r - j + 1)K\sigma^2\right\}\right) = 1 - F_{\chi^2(r-j+1)}\left((r - j + 1)K\right). \tag{6.5}$$

For all $j \in \{1, \dots, D_{m^*}\}$,

$$\begin{aligned}
 & \mathbb{P}\left(\left\{c_r + \dots + c_j > (r-j+1)K\sigma^2\right\}\right) \\
 &= \mathbb{P}\left(\left\{c_r + \dots + c_{D_{m^*+1}} + c_{D_{m^*}} + \dots + c_j > (r-j+1)K\sigma^2\right\}\right) \\
 &= \mathbb{P}\left(\left\{c_r + \dots + c_{D_{m^*+1}} + (\langle X\beta^*, u_{D_{m^*}} \rangle + \langle \varepsilon, u_{D_{m^*}} \rangle)^2 + \dots \right. \right. \\
 &\quad \left. \left. + (\langle X\beta^*, u_j \rangle + \langle \varepsilon, u_j \rangle)^2 > (r-j+1)K\sigma^2\right\}\right) \\
 &\stackrel{(**)}{\leq} \mathbb{P}\left(\left\{c_r + \dots + c_{D_{m^*+1}} + 2\langle X\beta^*, u_{D_{m^*}} \rangle^2 + 2\langle \varepsilon, u_{D_{m^*}} \rangle^2 + \dots \right. \right. \\
 &\quad \left. \left. + 2\langle X\beta^*, u_j \rangle^2 + 2\langle \varepsilon, u_j \rangle^2 > (r-j+1)K\sigma^2\right\}\right) \\
 &\leq \mathbb{P}\left(\left\{2c_r + \dots + 2c_{D_{m^*+1}} + 2\langle \varepsilon, u_{D_{m^*}} \rangle^2 + \dots + 2\langle \varepsilon, u_j \rangle^2 > (r-j+1)K\sigma^2 \right. \right. \\
 &\quad \left. \left. - 2\langle X\beta^*, u_{D_{m^*}} \rangle^2 - \dots - 2\langle X\beta^*, u_j \rangle^2\right\}\right) \\
 &\stackrel{(***)}{=} \mathbb{P}\left(\left\{2\sigma^2 Z_r^2 + \dots + 2\sigma^2 Z_{D_{m^*+1}}^2 + 2\sigma^2 Z_{D_{m^*}}^2 + \dots + 2\sigma^2 Z_j^2 \right. \right. \\
 &\quad \left. \left. > (r-j+1)K\sigma^2 - 2\langle X\beta^*, u_{D_{m^*}} \rangle^2 - \dots - 2\langle X\beta^*, u_j \rangle^2\right\}\right), \\
 &\quad \text{where } (Z_\ell)_{\ell \in \{j, \dots, r\}} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1). \\
 &= \mathbb{P}\left(\left\{Z_r^2 + \dots + Z_{D_{m^*+1}}^2 + Z_{D_{m^*}}^2 + \dots + Z_j^2 \right. \right. \\
 &\quad \left. \left. > \frac{(r-j+1)K}{2} - \frac{\langle X\beta^*, u_{D_{m^*}} \rangle^2}{\sigma^2} - \dots - \frac{\langle X\beta^*, u_j \rangle^2}{\sigma^2}\right\}\right) \\
 &= \mathbb{P}\left(\left\{X > \frac{(r-j+1)K}{2} - \frac{\langle X\beta^*, u_{D_{m^*}} \rangle^2}{\sigma^2} - \dots - \frac{\langle X\beta^*, u_j \rangle^2}{\sigma^2}\right\}\right), \\
 &\quad \text{for } X \sim \chi^2(r-j+1) \\
 &= 1 - F_{\chi^2(r-j+1)}\left(\frac{(r-j+1)K}{2} - \frac{\langle X\beta^*, u_{D_{m^*}} \rangle^2}{\sigma^2} - \dots - \frac{\langle X\beta^*, u_j \rangle^2}{\sigma^2}\right). \tag{6.6}
 \end{aligned}$$

(**) provides from $(a+b)^2 \leq 2(a^2+b^2)$, $\forall(a, b) \in \mathbb{R}$ and (***) is true since $\langle \varepsilon, u_i \rangle_{i \in \{1, \dots, r\}} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$.

So, from (6.4), (6.5) and (6.6), we deduce that for all $K > 0$ and for each $r \in \{D_{m^*} + 1, \dots, q\}$:

$$\begin{aligned}
 Q_r(K, \beta^*, \sigma^2) \leq & \min \left(1 - F_{\chi^2(1)}(K), \dots, 1 - F_{\chi^2(r-D_{m^*})}((r-D_{m^*})K), \right. \\
 & 1 - F_{\chi^2(r-D_{m^*+1})}\left(\frac{(r-D_{m^*}+1)K}{2} - \frac{\langle X\beta^*, u_{D_{m^*}} \rangle^2}{\sigma^2}\right), \\
 & 1 - F_{\chi^2(r-D_{m^*+2})}\left(\frac{(r-D_{m^*}+2)K}{2} - \frac{\langle X\beta^*, u_{D_{m^*}} \rangle^2}{\sigma^2} - \frac{\langle X\beta^*, u_{D_{m^*-1}} \rangle^2}{\sigma^2}\right), \\
 & \dots, \\
 & \left. 1 - F_{\chi^2(r)}\left(\frac{rK}{2} - \frac{\langle X\beta^*, u_{D_{m^*}} \rangle^2}{\sigma^2} - \frac{\langle X\beta^*, u_{D_{m^*-1}} \rangle^2}{\sigma^2} - \dots - \frac{\langle X\beta^*, u_1 \rangle^2}{\sigma^2}\right)\right).
 \end{aligned}$$

Trade-off between prediction and FDR for variable selection

Hence, an upper bound on $Q_r(K, \beta^*, \sigma^2)$ is obtained for all $K > 0$:

$$Q_r(K, \beta^*, \sigma^2) \leq \bar{f}_r(K, \beta^*, \sigma^2) \quad (6.7)$$

with:

$$\begin{aligned} \bar{f}_r(K, \beta^*, \sigma^2) = 1 - \max & \left(\max_{\ell \in \{1, \dots, r-D_{m^*}\}} \left(F_{\chi^2(\ell)}(\ell K) \right), \right. \\ & \left. \max_{\ell \in \{r-D_{m^*}+1, \dots, r\}} \left(F_{\chi^2(\ell)} \left(\frac{\ell K}{2} - \sum_{k=r-\ell+1}^{D_{m^*}} \frac{\langle X\beta^*, u_k \rangle^2}{\sigma^2} \right) \right) \right). \end{aligned} \quad (6.8)$$

- bounds on the FDR.

By combining (3.2), (6.2), (6.3), (6.7), (6.8) and (3.3), we obtain:

$$\sum_{r=D_{m^*}+1}^q \left(\frac{r-D_{m^*}}{r} P_r(K) \underline{f}_r(K, \beta^*, \sigma^2) \right) \leq \text{FDR}(\hat{m}(K))$$

and

$$\text{FDR}(\hat{m}(K)) \leq \sum_{r=D_{m^*}+1}^q \left(\frac{r-D_{m^*}}{r} P_r(K) \bar{f}_r(K, \beta^*, \sigma^2) \right),$$

which allows us to obtain Theorem 3.2 with $\forall K > 0$,

$$b(K, \beta^*, \sigma^2) = \sum_{r=D_{m^*}+1}^q \left(\frac{r-D_{m^*}}{r} P_r(K) \underline{f}_r(K, \beta^*, \sigma^2) \right)$$

and

$$B(K, \beta^*, \sigma^2) = \sum_{r=D_{m^*}+1}^q \left(\frac{r-D_{m^*}}{r} P_r(K) \bar{f}_r(K, \beta^*, \sigma^2) \right).$$

□

6.3 Strictly positive FDR.

Proof of Corollary 3.3.

From Theorem 3.2, we have $\forall K > 0$,

$$\text{FDR}(\hat{m}(K)) \geq \sum_{r=D_{m^*}+1}^q \left(\frac{r-D_{m^*}}{r} P_r(K) \underline{f}_r(K, \beta^*, \sigma^2) \right). \quad (6.9)$$

For the rest of the proof, we use the following Lemma:

Lemma 6.4 (Frank R. Kschischang [38]). *The complementary error function, $\text{erfc}(x)$, is defined, for $x \geq 0$, as:*

$$\text{erfc}(x) = 2 \left(1 - F_{\mathcal{N}(0, \frac{1}{2})}(x) \right)$$

where $F_{\mathcal{N}(0, \frac{1}{2})}$ designs the cumulative function of the centered Gaussian with the variance equals $\frac{1}{2}$.

Then,

$$\forall x \geq 0, \quad \frac{2e^{-x^2}}{\sqrt{\pi}(x + \sqrt{x^2 + 2})} \leq \text{erfc}(x) \leq \frac{e^{-x^2}}{\sqrt{\pi}x}.$$

We remark that for all $x \geq 0$, $1 - \Phi(x) = \frac{1}{2}\text{erfc}\left(\frac{x}{\sqrt{2}}\right)$. Then, for each $r \in \{D_{m^*} + 1, \dots, q\}$,

$$\begin{aligned}
 \underline{f}_r(K, \beta^*, \sigma^2) &= G_r + H_r \left(G_{r-1} + H_{r-1} (G_{r-2} + \dots + H_2 G_1) \right) \\
 &\geq G_r \\
 &= 2 \left(1 - \Phi(\sqrt{rK}) \right) \\
 &= \text{erfc}\left(\sqrt{\frac{rK}{2}}\right) \\
 &\stackrel{(**)}{\geq} \frac{2}{\sqrt{\pi} \left(\sqrt{\frac{rK}{2}} + \sqrt{\frac{rK}{2} + 2} \right)} e^{-\frac{rK}{2}} \\
 &= \frac{2\sqrt{2}}{\sqrt{\pi} \left(\sqrt{rK} + \sqrt{rK + 4} \right)} e^{-\frac{rK}{2}}. \tag{6.10}
 \end{aligned}$$

(**) is provided by Lemma 6.4. So, from (6.9) and (6.10), we obtain:

$$\forall K > 0, \text{FDR}(\hat{m}(K)) \geq \sum_{r=D_{m^*}+1}^q \left(\frac{r - D_{m^*}}{r} P_r(K) \frac{2\sqrt{2}}{\sqrt{\pi} \left(\sqrt{rK} + \sqrt{rK + 4} \right)} e^{-\frac{rK}{2}} \right).$$

This lower bound is strictly positive and since the $P_r(K)$ terms are all strictly positive too, we deduce that the FDR function is a strictly positive function. \square

6.4 Asymptotic analysis.

Proof of Corollary 3.4.

For all $r \in \{D_{m^*} + 1, \dots, q\}$ and by using the definitions from Theorem 3.2, for $\ell \in \{1, \dots, D_{m^*}\}$:

$$\begin{aligned}
 G_\ell &= 2 - \left(\Phi\left(\sqrt{\ell K} - \frac{\langle X\beta^*, u_\ell \rangle}{\sigma}\right) + \Phi\left(\sqrt{\ell K} + \frac{\langle X\beta^*, u_\ell \rangle}{\sigma}\right) \right) \\
 &\xrightarrow{K \rightarrow +\infty} 0;
 \end{aligned}$$

for $\ell \in \{2, \dots, D_{m^*}\}$:

$$\begin{aligned}
 H_\ell &= \Phi\left(\sqrt{\ell K} - \frac{\langle X\beta^*, u_\ell \rangle}{\sigma}\right) + \Phi\left(\sqrt{\ell K} + \frac{\langle X\beta^*, u_\ell \rangle}{\sigma}\right) - \\
 &\quad \left(\Phi\left(\sqrt{K} - \frac{\langle X\beta^*, u_\ell \rangle}{\sigma}\right) + \Phi\left(\sqrt{K} + \frac{\langle X\beta^*, u_\ell \rangle}{\sigma}\right) \right) \\
 &\xrightarrow{K \rightarrow +\infty} 0;
 \end{aligned}$$

and for $\ell \in \{D_{m^*} + 1, \dots, r\}$:

$$\begin{aligned}
 G_\ell &= 2 \left(1 - \Phi(\sqrt{\ell K}) \right) \xrightarrow{K \rightarrow +\infty} 0, \\
 H_\ell &= 2 \left(\Phi(\sqrt{\ell K}) - \Phi(\sqrt{K}) \right) \xrightarrow{K \rightarrow +\infty} 0;
 \end{aligned}$$

which provides that $\underline{f}_r(K, \beta^*, \sigma^2) \xrightarrow{K \rightarrow +\infty} 0$.

Moreover, $\bar{f}_r(K, \beta^*, \sigma^2) \xrightarrow{K \rightarrow +\infty} 0$. So, $Q_r(K, \beta^*, \sigma^2) \xrightarrow{K \rightarrow +\infty} 0$. In the same way, $P_r(K) \xrightarrow{K \rightarrow +\infty} 1$.

So, $P_r(K)Q_r(K, \beta^*, \sigma^2) \xrightarrow{K \rightarrow +\infty} 0$.

Finally, for each $r \in \{D_{m^*} + 1, \dots, q\}$, we deduce from (3.2) that

$$\text{FDR}(\hat{m}(K)) \xrightarrow{K \rightarrow +\infty} 0.$$

As for each $r \in \{D_{m^*} + 1, \dots, q\}$ $P_r(K) \xrightarrow{K \rightarrow +\infty} 1$, we deduce that for all $C_1 \in]0, 1[$, there exists $\tilde{L}_{C_1} > 0$ such that $\forall K > \tilde{L}_{C_1}$ and $\forall r \in \{D_{m^*} + 1, \dots, q\}$, we have $C_1 \leq P_r(K)$. For the following, we fix $C_1 \in]0, 1[$.

By using (6.2), (6.7) and $P_r(K) \leq 1$ for each $r \in \{D_{m^*} + 1, \dots, q\}$, we deduce that:

$$\forall K > \tilde{L}_{C_1}, \quad \text{FDR}(\hat{m}(K)) \geq C_1 \sum_{r=D_{m^*}+1}^q \left(\frac{r - D_{m^*}}{r} \underline{f}_r(K, \beta^*, \sigma^2) \right) \quad (6.11)$$

and

$$\forall K > 0, \quad \text{FDR}(\hat{m}(K)) \leq \sum_{r=D_{m^*}+1}^q \left(\frac{r - D_{m^*}}{r} \bar{f}_r(K, \beta^*, \sigma^2) \right). \quad (6.12)$$

- **Upper bound on \bar{f}_r :**

For each $r \in \{D_{m^*} + 1, \dots, q\}$ and for all $K > 0$:

$$\begin{aligned} \bar{f}_r(K, \beta^*, \sigma^2) &= 1 - \max \left(\max_{\ell \in \{1, \dots, r - D_{m^*}\}} \left(F_{\chi^2(\ell)}(\ell K) \right), \right. \\ &\quad \left. \max_{\ell \in \{r - D_{m^*} + 1, \dots, r\}} \left(F_{\chi^2(\ell)} \left(\frac{\ell K}{2} - \sum_{k=r-\ell+1}^{D_{m^*}} \frac{\langle X\beta^*, u_k \rangle^2}{\sigma^2} \right) \right) \right) \\ &= \min \left(\min_{\ell \in \{1, \dots, r - D_{m^*}\}} \left(\mathbb{P}(X_\ell > \ell K) \right), \right. \\ &\quad \left. \min_{\ell \in \{r - D_{m^*} + 1, \dots, r\}} \left(\mathbb{P} \left(Y_\ell > \frac{\ell K}{2} - \sum_{k=r-\ell+1}^{D_{m^*}} \frac{\langle X\beta^*, u_k \rangle^2}{\sigma^2} \right) \right), \right. \\ &\quad \left. \text{with } X_\ell \sim \chi^2(\ell) \text{ and } Y_\ell \sim \chi^2(\ell) \right) \\ &= \min \left(\min_{\ell \in \{1, \dots, r - D_{m^*}\}} \left(\mathbb{P}(X_\ell - \ell > \ell K - \ell) \right), \right. \\ &\quad \left. \min_{\ell \in \{r - D_{m^*} + 1, \dots, r\}} \left(\mathbb{P} \left(Y_\ell - \ell > \frac{\ell(K-2)}{2} - \sum_{k=r-\ell+1}^{D_{m^*}} \frac{\langle X\beta^*, u_k \rangle^2}{\sigma^2} \right) \right), \right. \\ &\quad \left. \text{with } X_\ell \sim \chi^2(\ell) \text{ and } Y_\ell \sim \chi^2(\ell). \right) \end{aligned} \quad (6.13)$$

So, for each $r \in \{D_{m^*} + 1, \dots, q\}$ and for all $K > 0$:

$$\bar{f}_r(K, \beta^*, \sigma^2) \leq \min_{\ell \in \{1, \dots, r - D_{m^*}\}} \left(\mathbb{P}(X_\ell - \ell > \ell K - \ell) \right), \quad \text{with } X_\ell \sim \chi^2(\ell).$$

By the exponential inequality of [39] for $X \sim \chi^2(\ell)$ and $\ell \in \mathbb{N}^*$:

$$\forall x \geq 0, \quad \mathbb{P}(X - \ell > 2\sqrt{\ell x} + 2x) \leq e^{-x}. \quad (6.14)$$

We apply (6.14) for each $\ell = 1, \dots, (r - D_{m^*})$ with $x = \frac{\ell}{4} \left(1 - \sqrt{2K - 1} \right)^2$ which is one solution of $2\sqrt{\ell x} + 2x = \ell K - \ell$ when $K > 1$. We obtain for all $K > 1$:

$$\begin{aligned} \min_{\ell \in \{1, \dots, r - D_{m^*}\}} \left(\mathbb{P}(X_\ell - \ell > \ell K - \ell) \right) &\leq \min_{\ell=1, \dots, (r - D_{m^*})} \left(e^{-\frac{\ell}{4} \left(1 - \sqrt{2K - 1} \right)^2} \right) \\ &\leq e^{-\frac{(r - D_{m^*}) \sqrt{2K - 1}}{2}} e^{-\frac{(r - D_{m^*}) K}{2}}. \end{aligned} \quad (6.15)$$

Trade-off between prediction and FDR for variable selection

So, from (6.12) and (6.15), we obtain for each $r \in \{D_{m^*} + 1, \dots, q\}$ and for all $K > 1$:

$$\begin{aligned} \text{FDR}(\hat{m}(K)) &\leq \sum_{r=D_{m^*}+1}^q \left(\frac{r - D_{m^*}}{r} e^{\frac{(r-D_{m^*})\sqrt{2K-1}}{2}} e^{-\frac{(r-D_{m^*})K}{2}} \right) \\ &\leq e^{-\frac{K}{2}} \sum_{r=D_{m^*}+1}^q \left(\frac{r - D_{m^*}}{r} e^{\frac{(r-D_{m^*})\sqrt{2K-1}}{2}} \right). \end{aligned}$$

For all $\eta > 0$ and $r \in \{D_{m^*} + 1, \dots, q\}$, $e^{\frac{(r-D_{m^*})\sqrt{2K-1}}{2}} = \underset{K \rightarrow +\infty}{o} (e^{\eta K})$.

Hence, $\forall \eta > 0$

$$\text{FDR}(\hat{m}(K)) = \underset{K \rightarrow +\infty}{o} \left(e^{-K(\frac{1}{2}-\eta)} \right),$$

which allows to obtain (3.6).

Proof of Remark 3.6:

The inequalities (6.11) and (6.12) are also true when $K \rightarrow +\infty$ and $\sigma \rightarrow 0$ with $\frac{1}{\sigma} = \underset{\sigma \rightarrow 0}{o} (\sqrt{K})$. To obtain the finest asymptotic upper bound (3.9), we start from the equation (6.13) and we consider the second term. Similar to previously, we apply (6.14) for each $\ell = r - D_{m^*} + 1, \dots, r$ with

$$x = \frac{\ell}{4} \left(1 - \sqrt{K-1 - \frac{2}{\ell} \sum_{k=r-\ell+1}^{D_{m^*}} \frac{\langle X\beta^*, u_k \rangle^2}{\sigma^2}} \right)^2,$$

which is one solution of

$$2\sqrt{\ell x} + 2x = \frac{\ell(K-2)}{2} - \sum_{k=r-\ell+1}^{D_{m^*}} \frac{\langle X\beta^*, u_k \rangle^2}{\sigma^2}$$

when $\sigma^2(K-1) > \frac{2}{r-D_{m^*}+1} \sum_{k=1}^{D_{m^*}} \langle X\beta^*, u_k \rangle^2 + 2$. This condition is valid since $\sigma \rightarrow 0$ with $\frac{1}{\sigma} = \underset{\sigma \rightarrow 0}{o} (\sqrt{K})$ leading to $\frac{1}{\sigma^2} = \underset{\sigma \rightarrow 0}{o} (K)$ and so $\sigma^2(K-1) \rightarrow +\infty$ when $K \rightarrow +\infty$. We obtain for all $K > 0$ such that

$$\sigma^2(K-1) > \frac{2}{r-D_{m^*}+1} \sum_{k=1}^{D_{m^*}} \langle X\beta^*, u_k \rangle^2 + 2:$$

$$\begin{aligned} &\min_{\ell \in \{r-D_{m^*}+1, \dots, r\}} \left(\mathbb{P} \left(Y_\ell - \ell > \frac{\ell(K-2)}{2} - \sum_{k=r-\ell+1}^{D_{m^*}} \frac{\langle X\beta^*, u_k \rangle^2}{\sigma^2} \right) \right) \\ &\leq \min_{\ell \in \{r-D_{m^*}+1, \dots, r\}} \left(e^{-\frac{\ell}{4} \left(1 - \sqrt{K-1 - \frac{2}{\ell} \sum_{k=r-\ell+1}^{D_{m^*}} \frac{\langle X\beta^*, u_k \rangle^2}{\sigma^2}} \right)^2} \right) \\ &\leq e^{\frac{1}{2} \sum_{k=1}^{D_{m^*}} \frac{\langle X\beta^*, u_k \rangle^2}{\sigma^2}} e^{\frac{\ell}{2} \sqrt{K-1 - \frac{2}{\ell} \sum_{k=1}^{D_{m^*}} \frac{\langle X\beta^*, u_k \rangle^2}{\sigma^2}}} e^{-\frac{\ell K}{4}}. \end{aligned} \tag{6.16}$$

(*) come from the fact that a minimum into a set is smaller than any value in the set. We choose the value corresponding for $\ell = 0$.

Trade-off between prediction and FDR for variable selection

So, from (6.12), (6.15) and (6.16), we obtain for each $r \in \{D_{m^*} + 1, \dots, q\}$ and for all $K > 1$ respecting

$$\sigma^2(K-1) > \frac{2}{r-D_{m^*}+1} \sum_{k=1}^{D_{m^*}} \langle X\beta^*, u_k \rangle^2 + 2:$$

$$\begin{aligned} \text{FDR}(\hat{m}(K)) &\leq \sum_{r=D_{m^*}+1}^q \left(\frac{r-D_{m^*}}{r} \min \left(e^{\frac{(r-D_{m^*})\sqrt{2K-1}}{2}} e^{-\frac{(r-D_{m^*})K}{2}}, \right. \right. \\ &\quad \left. \left. e^{\frac{1}{2} \sum_{k=1}^{D_{m^*}} \frac{\langle X\beta^*, u_k \rangle^2}{\sigma^2}} e^{\frac{r}{2} \sqrt{K-1-\frac{2}{r} \sum_{k=1}^{D_{m^*}} \frac{\langle X\beta^*, u_k \rangle^2}{\sigma^2}}} e^{-\frac{rK}{4}} \right) \right) \\ &= \min \left(\sum_{r=D_{m^*}+1}^q \left(\frac{r-D_{m^*}}{r} e^{\frac{(r-D_{m^*})\sqrt{2K-1}}{2}} e^{-\frac{(r-D_{m^*})K}{2}}, \right. \right. \\ &\quad \left. \left. \sum_{r=D_{m^*}+1}^q \left(\frac{r-D_{m^*}}{r} e^{\frac{1}{2} \sum_{k=1}^{D_{m^*}} \frac{\langle X\beta^*, u_k \rangle^2}{\sigma^2}} e^{\frac{r}{2} \sqrt{K-1-\frac{2}{r} \sum_{k=1}^{D_{m^*}} \frac{\langle X\beta^*, u_k \rangle^2}{\sigma^2}}} e^{-\frac{rK}{4}} \right) \right) \right) \\ &\leq \min \left(e^{-\frac{K}{2}} \sum_{r=D_{m^*}+1}^q \left(\frac{r-D_{m^*}}{r} e^{\frac{(r-D_{m^*})\sqrt{2K-1}}{2}}, \right. \right. \\ &\quad \left. \left. \sum_{r=D_{m^*}+1}^q \left(\frac{r-D_{m^*}}{r} e^{\frac{r}{2} \sqrt{K-1-\frac{2}{r} \sum_{k=1}^{D_{m^*}} \frac{\langle X\beta^*, u_k \rangle^2}{\sigma^2}}} e^{-\left(\frac{(D_{m^*}+1)K}{4} - \frac{1}{2\sigma^2} \sum_{k=1}^{D_{m^*}} \langle X\beta^*, u_k \rangle^2\right)} \right) \right). \end{aligned} \quad (6.17)$$

For all $\eta > 0$ and $r \in \{D_{m^*} + 1, \dots, q\}$, $e^{\frac{(r-D_{m^*})\sqrt{2K-1}}{2}} = o_{K \rightarrow +\infty} (e^{\eta K})$, independently of the value of σ^2 . Hence, the first term in (6.17) is $o(e^{-K(\frac{1}{2}-\eta)})$, $\forall \eta > 0$ when $K \rightarrow +\infty$ and $\sigma \rightarrow 0$ with $\frac{1}{\sigma} = o_{\sigma \rightarrow 0}(\sqrt{K})$.

For all $r \in \{D_{m^*} + 1, \dots, q\}$, $e^{\frac{r}{2} \sqrt{K-1-\frac{2}{r} \sum_{k=1}^{D_{m^*}} \frac{\langle X\beta^*, u_k \rangle^2}{\sigma^2}}} \leq e^{\frac{r}{2} \sqrt{K}}$. Moreover, for all $\tilde{\eta} > 0$ and $r \in \{D_{m^*} + 1, \dots, q\}$, $e^{\frac{r}{2} \sqrt{K}} = o_{K \rightarrow +\infty} (e^{\tilde{\eta} K})$, independently of the value of σ^2 . Hence, the second term in (6.17) is

$$o\left(e^{-\left(K \frac{(D_{m^*}+1-\tilde{\eta})}{4} - \frac{1}{2\sigma^2} \sum_{k=1}^{D_{m^*}} \langle X\beta^*, u_k \rangle^2\right)}\right), \forall \tilde{\eta} > 0 \text{ when } K \rightarrow +\infty \text{ and } \sigma \rightarrow 0 \text{ with } \frac{1}{\sigma} = o_{\sigma \rightarrow 0}(\sqrt{K}).$$

Hence,

$$\begin{aligned} \text{FDR}(\hat{m}(K)) &\leq \min \left(o\left(e^{-K(\frac{1}{2}-\eta)}\right), o\left(e^{-\left(K \frac{(D_{m^*}+1-\tilde{\eta})}{4} - \frac{1}{2\sigma^2} \sum_{k=1}^{D_{m^*}} \langle X\beta^*, u_k \rangle^2\right)}\right) \right) \\ &= o\left(e^{-\left(K \frac{(D_{m^*}+1-\tilde{\eta})}{4} - \frac{1}{2\sigma^2} \sum_{k=1}^{D_{m^*}} \langle X\beta^*, u_k \rangle^2\right)}\right). \end{aligned}$$

$\forall (\eta, \tilde{\eta}) > 0$ when $K \rightarrow +\infty$ and $\sigma \rightarrow 0$ with $\frac{1}{\sigma} = o_{\sigma \rightarrow 0}(\sqrt{K})$; which allows us to obtain (3.9).

- **Lower bound on \underline{f}_r :**

From (6.10) and (6.11), we obtain:

$$\begin{aligned}
 \forall K > \tilde{L}_{C_1}, \text{FDR}(\hat{m}(K)) &\geq C_1 \sum_{r=D_{m^*}+1}^q \left(\frac{r - D_{m^*}}{r} \frac{2\sqrt{2}}{\sqrt{\pi}(\sqrt{rK} + \sqrt{rK+4})} e^{-\frac{rK}{2}} \right) \\
 &\geq C_1 \frac{2\sqrt{2}}{\sqrt{\pi}(\sqrt{qK} + \sqrt{qK+4})} \frac{1}{D_{m^*} + 1} \sum_{r=D_{m^*}+1}^q \left(e^{-\frac{rK}{2}} \right) \\
 &\stackrel{(*)}{\geq} C_1 \frac{2\sqrt{2}}{\sqrt{\pi}(\sqrt{qK} + \sqrt{qK+4})} \frac{1}{D_{m^*} + 1} e^{-\frac{(D_{m^*}+1)K}{2}} \\
 &= \frac{2\sqrt{2}C_1}{\sqrt{\pi}(D_{m^*} + 1)} \frac{1}{\sqrt{qK} + \sqrt{qK+4}} e^{-K \frac{(D_{m^*}+1)}{2}}.
 \end{aligned}$$

(*) is true since each term in the sum is positive, so, the sum is larger than one of them.

For all $\eta > 0$, $\exists \tilde{C}_\eta > 0$, $\exists \tilde{L}_\eta > 0$ such that $\forall K > \tilde{L}_\eta$, we have $\tilde{C}_\eta e^{-\eta K} \leq \frac{1}{\sqrt{qK} + \sqrt{qK+4}}$.

So,

$$\begin{aligned}
 \forall \eta > 0, \exists \tilde{C}_\eta > 0, \exists \tilde{L}_\eta > 0, \forall K > \max(\tilde{L}_{C_1}, \tilde{L}_\eta), \\
 \text{FDR}(\hat{m}(K)) &\geq \frac{2\sqrt{2}C_1}{\sqrt{\pi}(D_{m^*} + 1)} \tilde{C}_\eta e^{-K \left(\frac{D_{m^*}+1+2\eta}{2} \right)},
 \end{aligned}$$

which gives (3.7) with $C_\eta = \frac{2\sqrt{2}C_1}{\sqrt{\pi}(D_{m^*}+1)} \tilde{C}_\eta$ and $L_\eta = \max(\tilde{L}_{C_1}, \tilde{L}_\eta)$.

Formula (3.8) automatically follows from (3.6) and (3.7). □

6.5 General bounds.

Proof of Corollary 3.7.

By taking $u_j = X_j$, $\forall j \in \{1, \dots, q\}$, then $(X_1, \dots, X_q, u_{q+1}, \dots, u_n)$ is an orthonormal basis of \mathbb{R}^n . Consequently, $\forall j \in \{1, \dots, q\}$, $\langle X\beta^*, u_j \rangle = \langle X\beta^*, X_j \rangle = \beta_j^*$, which concludes the proof. □

6.6 Estimation of the PR.

Proof of Proposition 4.1.

Let us observe that for all $K > 0$ and $K' > 0$:

$$\begin{aligned}
 &\mathbb{E}[\|Y - X\hat{\beta}_{\hat{m}(K)}\|_2^2] - \mathbb{E}[\|Y - X\hat{\beta}_{\hat{m}(K')}\|_2^2] \\
 &= \mathbb{E}[\|X\hat{\beta}_{\hat{m}(K)}\|_2^2] - \mathbb{E}[\|X\hat{\beta}_{\hat{m}(K')}\|_2^2] + 2\mathbb{E}[\langle Y, X\hat{\beta}_{\hat{m}(K')} - X\hat{\beta}_{\hat{m}(K)} \rangle] \\
 &= \mathbb{E}[\|X\hat{\beta}_{\hat{m}(2)} - X\hat{\beta}_{\hat{m}(K)}\|_2^2] - \mathbb{E}[\|X\hat{\beta}_{\hat{m}(2)} - X\hat{\beta}_{\hat{m}(K')}\|_2^2] \\
 &\quad + 2\mathbb{E}[\langle X\hat{\beta}_{\hat{m}(K)} - X\hat{\beta}_{\hat{m}(K')}, X\hat{\beta}_{\hat{m}(2)} \rangle] + 2\mathbb{E}[\langle Y, X\hat{\beta}_{\hat{m}(K')} - X\hat{\beta}_{\hat{m}(K)} \rangle] \\
 &= \mathbb{E}[\|X\hat{\beta}_{\hat{m}(2)} - X\hat{\beta}_{\hat{m}(K)}\|_2^2] - \mathbb{E}[\|X\hat{\beta}_{\hat{m}(2)} - X\hat{\beta}_{\hat{m}(K')}\|_2^2] \\
 &\quad - 2\mathbb{E}[\langle X\hat{\beta}_{\hat{m}(2)} - Y, X\hat{\beta}_{\hat{m}(K')} - X\hat{\beta}_{\hat{m}(K)} \rangle].
 \end{aligned} \tag{6.18}$$

The constant 2 allows to get the optimal asymptotic control of (2.3). Consequently, $\|Y - X\hat{\beta}_{\hat{m}(2)}\|_2$ is close to 0 and $X\hat{\beta}_{\hat{m}(2)} - Y$ almost belongs to the subspace $\text{Im}(X)^\perp$ since $X\hat{\beta}_{\hat{m}(2)}$ is close to $\Pi_{\text{Im}(X)}(Y)$. Since $X\hat{\beta}_{\hat{m}(K)}$ and $X\hat{\beta}_{\hat{m}(K')}$ belongs to $\text{Im}(X)$, the term $\mathbb{E}[\langle X\hat{\beta}_{\hat{m}(2)} - Y, X\hat{\beta}_{\hat{m}(K')} - X\hat{\beta}_{\hat{m}(K)} \rangle]$ in (6.18) is close to 0 and is negligible compared to the two others. So, the dynamics of $\mathbb{E}[\|X\hat{\beta}_{\hat{m}(2)} - X\hat{\beta}_{\hat{m}(K)}\|_2^2]$ with respect to positive K is close to the one of $\mathbb{E}[\|Y - X\hat{\beta}_{\hat{m}(K)}\|_2^2]$. □

7 Extensive simulation study to justify the observable estimations.

This section is a complement to Section 4 and presents an extensive simulation study.

7.1 Description of the simulation protocol

Description of the data simulation. Given values of n and p , we simulate $Y \sim \mathcal{N}(\beta^*, I_n)$ where β^* is a vector satisfying $\beta_j^* \leq \beta_{j+1}^*$ for all $j \in \{1, \dots, D_{m^*} - 1\}$ to get ordered active variables. We consider four scenarios, described in Table 1, where values of D_{m^*} , β^* , n and σ^2 vary and where the number of variables p is always equal to 50.

Scenario with $p = 50$	Active variable number	Non-zero coefficient amplitude in β^*	Observation number	Noise amplitude
(i) Sparsity	$D_{m^*} \in \{1, 10, 20\}$	$\beta_{D_{m^*}}^* = 2,$ $\forall j \in \{1, \dots, D_{m^*} - 1\}$ $\beta_j^* \sim \text{Unif}(\beta_{j+1}^* + 0.5, \beta_{j+1}^* + 1.5)$	$n = 50$	$\sigma^2 = 1$
(ii) Complexity	$D_{m^*} = 10$	$\beta_{10}^* = 2$ with $\forall j \in \{1, \dots, 9\}$ $\beta_j^* \sim \text{Unif}(\beta_{j+1}^* + 0.5, \beta_{j+1}^* + 1.5)$ $\beta_{10}^* = \frac{2}{10}$ with $\forall j \in \{1, \dots, 9\},$ $\beta_j^* \sim \text{Unif}(\beta_{j+1}^* + 0.05, \beta_{j+1}^* + 0.15)$ $\beta_{10}^* = 2$ with $\forall j \in \{1, \dots, 9\}$ $\beta_j^* \sim \text{Unif}(\beta_{j+1}^* + 0.05, \beta_{j+1}^* + 0.15)$	$n = 50$	$\sigma^2 = 1$
(iii) High-dimension	$D_{m^*} = 10$	$\beta_{D_{m^*}}^* = 2,$ $\forall j \in \{1, \dots, 9\}$ $\beta_j^* \sim \text{Unif}(\beta_{j+1}^* + 0.5, \beta_{j+1}^* + 1.5)$	$n \in \{30, 50, 300\}$	$\sigma^2 = 1$
(iv) Noise	$D_{m^*} = 10$	$\beta_{D_{m^*}}^* = 2,$ $\forall j \in \{1, \dots, 9\}$ $\beta_j^* \sim \text{Unif}(\beta_{j+1}^* + 0.5, \beta_{j+1}^* + 1.5)$	$n = 50$	$\sigma^2 \in \{0.1, 1, 4\}$

Table 1: Description of the four scenarios.

The scenario (i) allows us to evaluate the impact of the sparsity of the parameter β^* . The scenario (ii) allows us to evaluate how the values of the non-zero coefficients in β^* complicate the identification of the active variables. In particular, the non-zero coefficients are close and, in the second configuration, some of them are smaller than the noise level σ . The scenario (iii) allows us to evaluate the behavior of our method in a high-dimensional context through the variation of the number of observations n , either smaller, equal or larger than the number of variables p . The last scenario (iv) allows us to evaluate the impact of the noise amplitude through different values of σ^2 .

Note that for a fair comparison, the datasets where $n = 30$ in scenario (iii) are included in those where $n = 50$ which are included in those where $n = 300$. Moreover, for the sake of reproducibility, the seed of the random number generator is identically fixed for each scenario.

The toy data set. We call the *toy data set* the data set where $n = p = 50$, $D_{m^*} = 10$, $\beta_{10}^* = 2$ and $\forall j \in \{1, \dots, 9\}$, $\beta_j^* \sim \text{Unif}(\beta_{j+1}^* + 0.5, \beta_{j+1}^* + 1.5)$. It corresponds to the reference data set in all scenarios.

Empirical estimations. For the empirical estimations, we simulate \mathcal{D} a set of 1000 data sets for each scenario. For each $d \in \mathcal{D}$ and for all $K > 0$, the selected model $\hat{m}^d(K)$ is obtained from (Y^d, X^d) . Since m^* is known, the

quantity $\text{FDP}(\hat{m}^d(K))$ is calculable for each $d \in \mathcal{D}$ and the empirical estimator of $\text{FDR}(\hat{m}(K))$ is the average of the $\text{FDP}(\hat{m}^d(K))$. Concerning PR, we simulate $\tilde{\mathcal{D}}$ a new set of 1000 data sets for each scenario. New \tilde{Y}^d are generated on $\tilde{\mathcal{D}}$, from the model (1.1), and by using the X^d on \mathcal{D} to respect the fixed design assumption. The selected models $\hat{m}^d(K)$ and the $\hat{\beta}_{\hat{m}^d(K)}^d$ estimators are extracted by solving (2.2) from the training sets (Y^d, X^d) on \mathcal{D} . The PR is evaluated from the validation sets (\tilde{Y}^d, X^d) on $\tilde{\mathcal{D}}$ by the mean squared error:

$$\text{MSE}(\hat{m}^d(K)) = \frac{1}{n} \sum_{i=1}^n \left(\tilde{Y}_i^d - \sum_{j=1}^p x_{ij}^d \hat{\beta}_{\hat{m}^d(K)_j}^d \right)^2. \quad (7.1)$$

The empirical estimator of $\text{PR}(\hat{m}(K))$ is the average of the $\text{MSE}(\hat{m}^d(K))$.

To validate the quality of the empirical estimations, the central limit theorem is applied to get the 95% asymptotic confidence intervals:

$$\left[\text{FDR}(\hat{m}(K)) - 1.96 \frac{\hat{\sigma}}{\sqrt{1000}}, \text{FDR}(\hat{m}(K)) + 1.96 \frac{\hat{\sigma}}{\sqrt{1000}} \right]$$

and

$$\left[\text{PR}(\hat{m}(K)) - 1.96 \frac{\hat{\sigma}}{\sqrt{1000}}, \text{PR}(\hat{m}(K)) + 1.96 \frac{\hat{\sigma}}{\sqrt{1000}} \right],$$

where $\hat{\sigma}$ is the unbiased empirical estimator of the standard deviation σ . Since their width do not exceed 0.011 and 0.07 for respectively the FDR and the PR, they are tight, meaning that the empirical estimations are closed to the theoretical quantities $\text{FDR}(\hat{m}(K))$ and $\text{PR}(\hat{m}(K))$.

7.2 Estimation of the theoretical FDR

This subsection completes Subsection 4.1. We present the slope heuristic principle and an analyse of the $\hat{\sigma}^2$, obtained by the slope heuristics, is processed. Then, a large simulation study is performed to justify the choice of $\hat{\beta}_{\hat{m}(4)}$ to estimate β^* in the upper bound of the FDR.

The FDR bounds of Theorem 3.2 depend on the P_r , the $f_r(K, \beta^*, \sigma^2)$ and the $\bar{f}_r(K, \beta^*, \sigma^2)$ quantities. Concerning the P_r quantities, they do not depend on the data as soon as r is given. They can be estimated once and for all without any dataset. For each $1 \leq r \leq q$, P_r is estimated by generating 5000 independent standard Gaussian vectors $(Z_k)_{k \in \{r+1, \dots, q\}}$ and by counting for each vector the number of times that $Z_k^2 < K(\ell - r)$ for each $\ell \in \{r+1, \dots, q\}$.

Concerning the $f_r(K, \beta^*, \sigma^2)$ and $\bar{f}_r(K, \beta^*, \sigma^2)$ quantities, they depend on β^* and σ^2 , both unknown.

The slope heuristic to estimate σ^2 . The slope heuristic principle, introduced in [12], is that when D_m is large enough, the empirical least squares values $\frac{1}{n} \|Y - X \hat{\beta}_m\|_2^2$ are almost equal to $-\frac{1}{2n} K \sigma^2 D_m$ plus an additive constant independent of n and m . Hence, it is possible to estimate σ^2 from the dataset by the multiplicative coefficient of the affine behavior between the empirical least squares and $-\frac{K}{2n} D_m$ for D_m large enough. We use the function `capushe` of the R package `capushe` (version 1.1.1) [13] with parameters set to the default values.

Some substitutes of β^* . According to [12], $\hat{\beta}_{\hat{m}(K)}$ is a good estimator of β^* in a predictive point of view when K is equal or close to 2. We propose to test the estimators $\hat{\beta}_{\hat{m}(\tilde{K})}$ for $\tilde{K} \in \{1, 1.5, 2, 2.5, 3, 3.5, 4, 4.5, 5, \log(n)\}$ to replace β^* in the lower and upper bounds $b(K, \beta^*, \sigma^2)$ and $B(K, \beta^*, \sigma^2)$.

To determine the best constant \tilde{K} among $\{1, 1.5, 2, 2.5, 3, 3.5, 4, 4.5, 5, \log(n)\}$, we evaluate all $b(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ and $B(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ on the sets \mathcal{D} from the four scenarios described in Subsection 7.1. To take into account the randomness of $b(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ and $B(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$, the model collection generation and model selection given by (2.2) are processed on a new data set independent of \mathcal{D} for the four scenarios.

To evaluate the error by replacing $b(K, \beta^*, \sigma^2)$ and $B(K, \beta^*, \sigma^2)$ with their estimation $b(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ and $B(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$, we propose to evaluate the relative changes defined by: $\forall K > 0$,

$$\frac{b(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2) - b(K, \beta^*, \sigma^2)}{b(K, \beta^*, \sigma^2)}$$

for the lower bound and by:

$$\frac{B(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2) - B(K, \beta^*, \sigma^2)}{B(K, \beta^*, \sigma^2)}$$

for the upper bound. To ensure that $B(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ is above $B(K, \beta^*, \sigma^2)$ and so above the FDR, positive relative change values and as close to 0 as possible are expected. Concerning the lower bounds, negative relative change values are expected to ensure that $b(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ is below $B(K, \beta^*, \sigma^2)$ and so below the FDR. To take into account randomness of the $b(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ and $B(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ terms, we evaluate for all K the relative standard deviation, defined by the standard deviation divided by the mean, by calculated the variance of bounds $b(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ and $B(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ evaluated on 100 new data sets generated independently of \mathcal{D} . The relative standard deviation values are expected to be as close to 0 as possible.

Figures 4-9 are plotted from the *toy data set*. In Figures 4 and 5, the empirical estimation of the FDR ($\hat{m}(K)$) and the quantities $b(K, \beta^*, \sigma^2)$, $B(K, \beta^*, \sigma^2)$, $b(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ and $B(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ are plotted on a grid of positive K . Relative changes and relative standard deviations for the lower bounds $b(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ are plotted in Figures 6 and 7. Relative changes and relative standard deviations for the upper bounds $B(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ are plotted in Figures 8 and 9. The graphs of all others \mathcal{D} of the 4 scenarios described in Subsection 7.1 are provided in the supplementary material available in ¹.

The lower bounds: For $\tilde{K} > 1$, the relative change values are positive until achieving more than 2 for large K (Figure 6) and the estimated lower bounds curves can be larger than the theoretical one. The relative standard deviation functions increase quickly whatever the value of \tilde{K} suggesting that fluctuations around the mean are not negligible (Figure 7).

The upper bounds: For $\tilde{K} > 1$, the relative change functions are always positive and do not exceed 0.11 meaning that the $B(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ curves are close to $B(K, \beta^*, \sigma^2)$ for all $K > 0$ (Figure 8). For data sets \mathcal{D} other than the *toy data set* (Figures are available in Supplementary material ¹), the relative change values are always small but can be negative. However, it happens very rarely for $\tilde{K} \geq 4$ and in this case, values are low enough (smaller than -0.02%) to ensure that the empirical FDR estimation curves do not exceed the $B(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ terms. Concerning the relative standard deviation functions (Figures 9), the larger \tilde{K} , the smaller the values, except for the scenario (ii) with the third configuration where values increase after $\tilde{K} \geq 4.5$. For $\tilde{K} \geq 3.5$, the relative standard deviation values are around 0.2 for all the scenarios except for scenario (ii) with the second configuration (can achieve 0.8) and with the third configuration (can achieve 1). Thus, for a value of $\tilde{K} \in \{3.5, \log(n), 4, 4.5, 5\}$ and eventually except for the two extreme scenarios, fluctuations around the mean are small, meaning that the upper bound estimations are stable.

To conclude, we drop the lower bound to implement our data-driven algorithm for hyperparameter calibration since $b(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ functions can be larger than the theoretical FDR one. To control the FDR, only an upper bound control is sufficient. The best results for $B(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ are obtained with the hyperparameter $\tilde{K} = 4$, where the relative change values are almost always positive, small enough to guarantee that the $B(K, \hat{\beta}_{\hat{m}(4)}, \hat{\sigma}^2)$ are larger than the theoretical FDR, and the relative standard deviation values are the smallest ones whatever the scenarios. So, the estimator used in our algorithm to replace β^* in the upper bound of the FDR is $\hat{\beta}_{\hat{m}(4)}$. The value of the hyperparameter $\tilde{K} = 4$ is not surprising since the value of $D_{\hat{m}}$ has to be small enough in (3.5) to get an upper bound $B(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ larger than the theoretical upper one. So, the penalization function has to be large enough in (2.2).

Algorithm 1 gets $K = 2.7$ for scenario (i) with $D_m^* = 20$; $K = 4.8$ for scenario (i) with $D_m^* = 1$ and for scenario (ii) with the second configuration; and $K = 3.3$ for all the others. Thus, the hyperparameter K returned by the algorithm is strictly larger than the commonly used constant 2.

¹https://github.com/PerrineLacroix/Trade_off_FDR_PR

8 Acknowledgments

The authors warmly thank and are grateful to Pascal Massart (Laboratoire de Mathématiques d'Orsay, Université Paris-Saclay) for helpful discussions and valuable comments.

This research is supported in part by a public grant as part of the Investissement d'avenir project, reference ANR-11-LABX-0056-LMH, LabEx LMH. IPS2 benefits from the support of the LabEx Saclay Plant Sciences-SPS (ANR-17-EUR-0007).

9 Supplementary data

All the R scripts are available at https://github.com/PerrineLacroix/Trade_off_FDR_PR.

The graphs for the bounds $B(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ from the 4 scenarios described in Subsection 7.1 are provided in the supplementary material available in https://github.com/PerrineLacroix/Trade_off_FDR_PR. It is complementary to Subsection 7.2.

References

- [1] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- [2] Florentina Bunea, Alexandre B Tsybakov, and Marten H Wegkamp. Aggregation for gaussian regression. *The Annals of Statistics*, 35(4):1674–1697, 2007.
- [3] Florentina Bunea, Alexandre Tsybakov, and Marten Wegkamp. Sparsity oracle inequalities for the lasso. *Electronic Journal of Statistics*, 1:169–194, 2007.
- [4] N. Meinshausen and P. Bühlmann. Stability selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(4):417–473, 2010.
- [5] F. Bach. Bolasso: model consistent lasso estimation through the bootstrap. In *Proceedings of the 25th international conference on Machine learning*, pages 33–40, 2008.
- [6] Seymour Geisser. The predictive sample reuse method with applications. *Journal of the American statistical Association*, 70(350):320–328, 1975.
- [7] Sylvain Arlot and Alain Celisse. A survey of cross-validation procedures for model selection. *Statistics surveys*, 4:40–79, 2010.
- [8] H Akaike. Information theory and an extension of maximum likelihood principle. In *Proc. 2nd Int. Symp. on Information Theory*, pages 267–281, 1973.
- [9] Colin L Mallows. Some comments on cp. *Technometrics*, 42(1):87–94, 2000.
- [10] G. Schwarz. Estimating the dimension of a model. *The annals of statistics*, 6(2):461–464, 1978.
- [11] J. Chen and Z. Chen. Extended bayesian information criteria for model selection with large model spaces. *Biometrika*, 95(3):759–771, 2008.
- [12] L. Birgé and P. Massart. Minimal penalties for gaussian model selection. *Probability theory and related fields*, 138(1-2):33–73, 2007.
- [13] J.P. Baudry, C. Maugis, and B. Michel. Slope heuristics: overview and implementation. *Statistics and Computing*, 22(2):455–470, 2012.
- [14] Y. Baraud, C. Giraud, S. Huet, et al. Gaussian model selection with an unknown variance. *The Annals of Statistics*, 37(2):630–672, 2009.
- [15] C. Giraud, S. Huet, N. Verzelen, et al. High-dimensional regression with unknown variance. *Statistical Science*, 27(4):500–518, 2012.
- [16] C Bonferroni. Teoria statistica delle classi e calcolo delle probabilita. *Pubblicazioni del R Istituto Superiore di Scienze Economiche e Commerciali di Firenze*, 8:3–62, 1936.
- [17] R John Simes. An improved bonferroni procedure for multiple tests of significance. *Biometrika*, 73(3):751–754, 1986.
- [18] 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.
- [19] Yoav Benjamini and Daniel Yekutieli. The control of the false discovery rate in multiple testing under dependency. *Annals of statistics*, pages 1165–1188, 2001.
- [20] J. Storey, J. Taylor, and D. Siegmund. Strong control, conservative point estimation and simultaneous conservative consistency of false discovery rates: a unified approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 66(1):187–205, 2004.
- [21] J. Romano, A. Shaikh, and M. Wolf. Control of the false discovery rate under dependence using the bootstrap and subsampling. *Test*, 17(3):417–442, 2008.
- [22] D. Leung and W. Sun. Zap: z -value adaptive procedures for false discovery rate control with side information. *arXiv preprint arXiv:2108.12623*, 2021.
- [23] R. Barber, E. Candès, et al. Controlling the false discovery rate via knockoffs. *The Annals of Statistics*, 43(5):2055–2085, 2015.
- [24] Shuheng Zhou. Thresholding procedures for high dimensional variable selection and statistical estimation. *Advances in Neural Information Processing Systems*, 22:2304–2312, 2009.
- [25] Richard Berk, Lawrence Brown, Andreas Buja, Kai Zhang, and Linda Zhao. Valid post-selection inference. *The Annals of Statistics*, pages 802–837, 2013.

- [26] J. Lee, D. Sun, Y. Sun, and J. Taylor. Exact post-selection inference, with application to the lasso. *The Annals of Statistics*, 44(3):907–927, 2016.
- [27] S. Hyun, M. G’sell, and R. Tibshirani. Exact post-selection inference for the generalized lasso path. *Electronic Journal of Statistics*, 12(1):1053–1097, 2018.
- [28] Y. Chen, S. Jewell, and D. Witten. More powerful selective inference for the graph fused lasso, 2021.
- [29] V. Duy and I. Takeuchi. More powerful conditional selective inference for generalized lasso by parametric programming. *arXiv preprint arXiv:2105.04920*, 2021.
- [30] Dongliang Zhang, Abbas Khalili, and Masoud Asgharian. Post-model-selection inference in linear regression models: An integrated review. *Statistics Surveys*, 16:86–136, 2022.
- [31] C. Genovese and L. Wasserman. Operating characteristics and extensions of the false discovery rate procedure. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 64(3):499–517, 2002.
- [32] C. Genovese, L. Wasserman, et al. A stochastic process approach to false discovery control. *The Annals of Statistics*, 32(3):1035–1061, 2004.
- [33] F. Abramovich, Y. Benjamini, D. Donoho, I. Johnstone, et al. Adapting to unknown sparsity by controlling the false discovery rate. *The Annals of Statistics*, 34(2):584–653, 2006.
- [34] M. Bogdan, E. Berg, W. Su, and E. Candès. Statistical estimation and testing via the sorted l_1 norm. *arXiv preprint arXiv:1310.1969*, 2013.
- [35] Michał Kos and Małgorzata Bogdan. On the asymptotic properties of slope. *Sankhya A*, 82(2):499–532, 2020.
- [36] L. Birgé and P. Massart. Gaussian model selection. *Journal of the European Mathematical Society*, 3(3):203–268, 2001.
- [37] Y. Yang. Can the strengths of aic and bic be shared? a conflict between model identification and regression estimation. *Biometrika*, 92(4):937–950, 2005.
- [38] Frank R Kschischang. The complementary error function. *Online, April*, 2017.
- [39] Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *Annals of Statistics*, pages 1302–1338, 2000.

Trade-off between prediction and FDR for variable selection

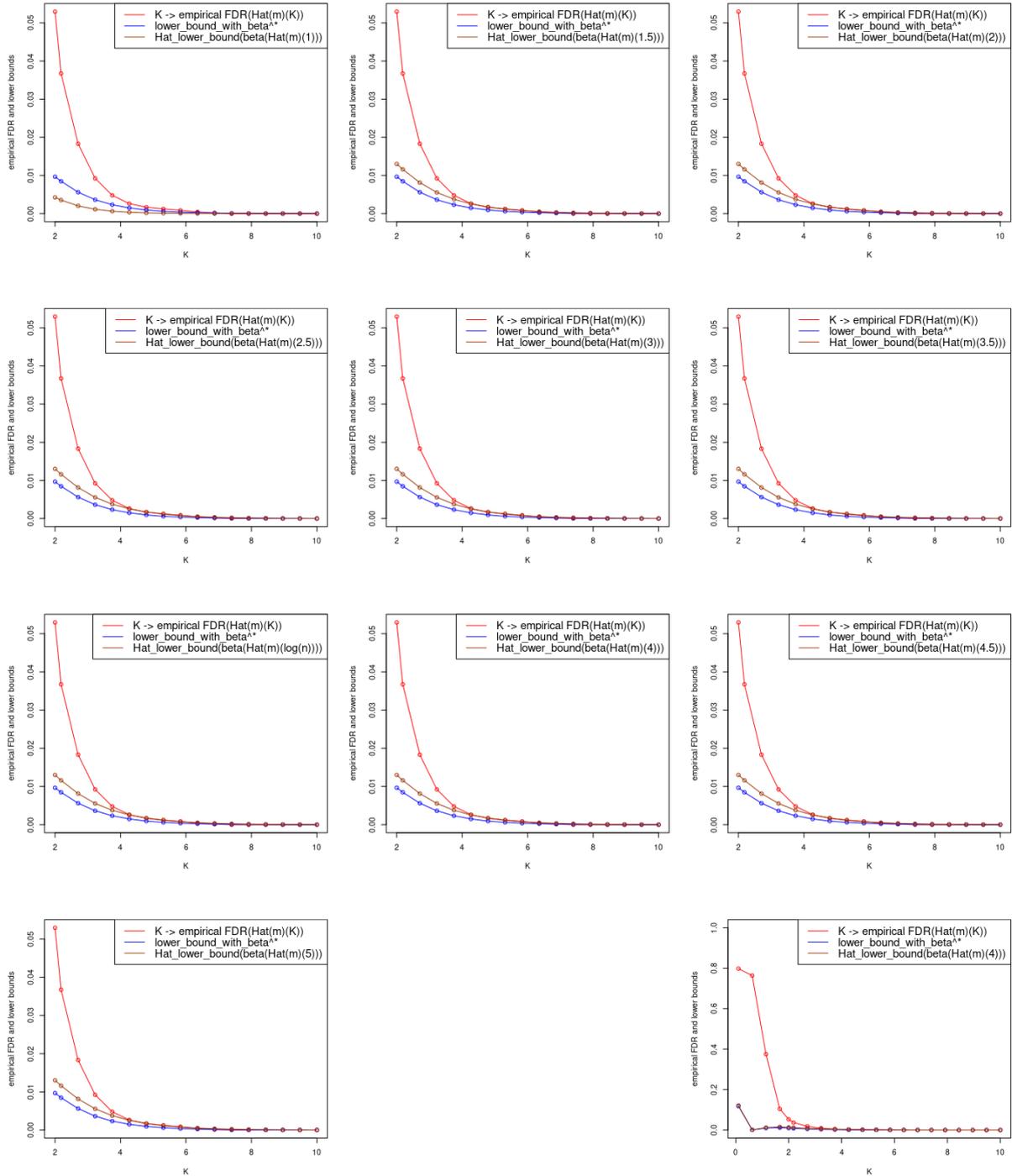


Figure 4: Comparison of the empirical estimation of the FDR, the function $b(K, \beta^*, \sigma^2)$ under the orthogonal design matrix X and the function $b(K, \hat{\beta}_{\hat{m}(\bar{K})}, \hat{\sigma}^2)$ with respectively $\hat{\beta}_{\hat{m}(1)}$, $\hat{\beta}_{\hat{m}(1.5)}$, $\hat{\beta}_{\hat{m}(2)}$, $\hat{\beta}_{\hat{m}(2.5)}$, $\hat{\beta}_{\hat{m}(3)}$, $\hat{\beta}_{\hat{m}(3.5)}$, $\hat{\beta}_{\hat{m}(4)}$, $\hat{\beta}_{\hat{m}(4.5)}$, $\hat{\beta}_{\hat{m}(5)}$ and $\hat{\beta}_{\hat{m}(\log(n))}$. The terms $b(K, \hat{\beta}_{\hat{m}(\bar{K})}, \hat{\sigma}^2)$ are calculating from only one dataset, independent of those used for the empirical estimations. For a better readability, we plot curves only for $K \geq 2$; but at the bottom right is the entire curve for $\bar{K} = 4$.

Trade-off between prediction and FDR for variable selection

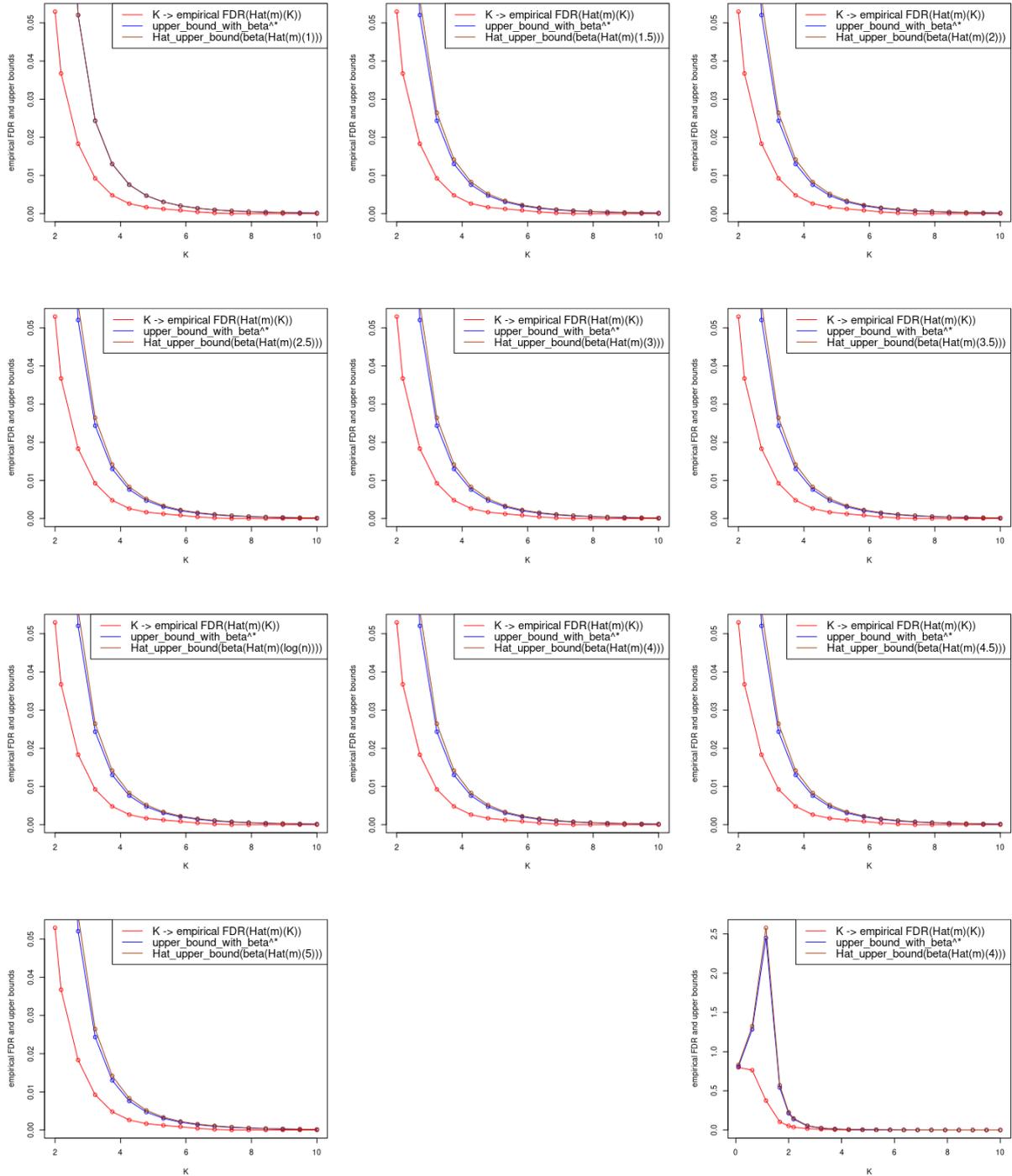


Figure 5: Comparison of the empirical estimation of the FDR, the function $B(K, \beta^*, \sigma^2)$ under the orthogonal design matrix X and the function $B(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ with respectively $\hat{\beta}_{\hat{m}(1)}, \hat{\beta}_{\hat{m}(1.5)}, \hat{\beta}_{\hat{m}(2)}, \hat{\beta}_{\hat{m}(2.5)}, \hat{\beta}_{\hat{m}(3)}, \hat{\beta}_{\hat{m}(3.5)}, \hat{\beta}_{\hat{m}(4)}, \hat{\beta}_{\hat{m}(4.5)}, \hat{\beta}_{\hat{m}(5)}$ and $\hat{\beta}_{\hat{m}(\log(n))}$. The terms $B(K, \hat{\beta}_{\hat{m}(\tilde{K})}, \hat{\sigma}^2)$ are calculating from only one dataset independent of those used for the empirical estimations. For a better readability, we plot curves only for $K \geq 2$; but at the bottom right is the entire curve for $\tilde{K} = 4$.

Trade-off between prediction and FDR for variable selection

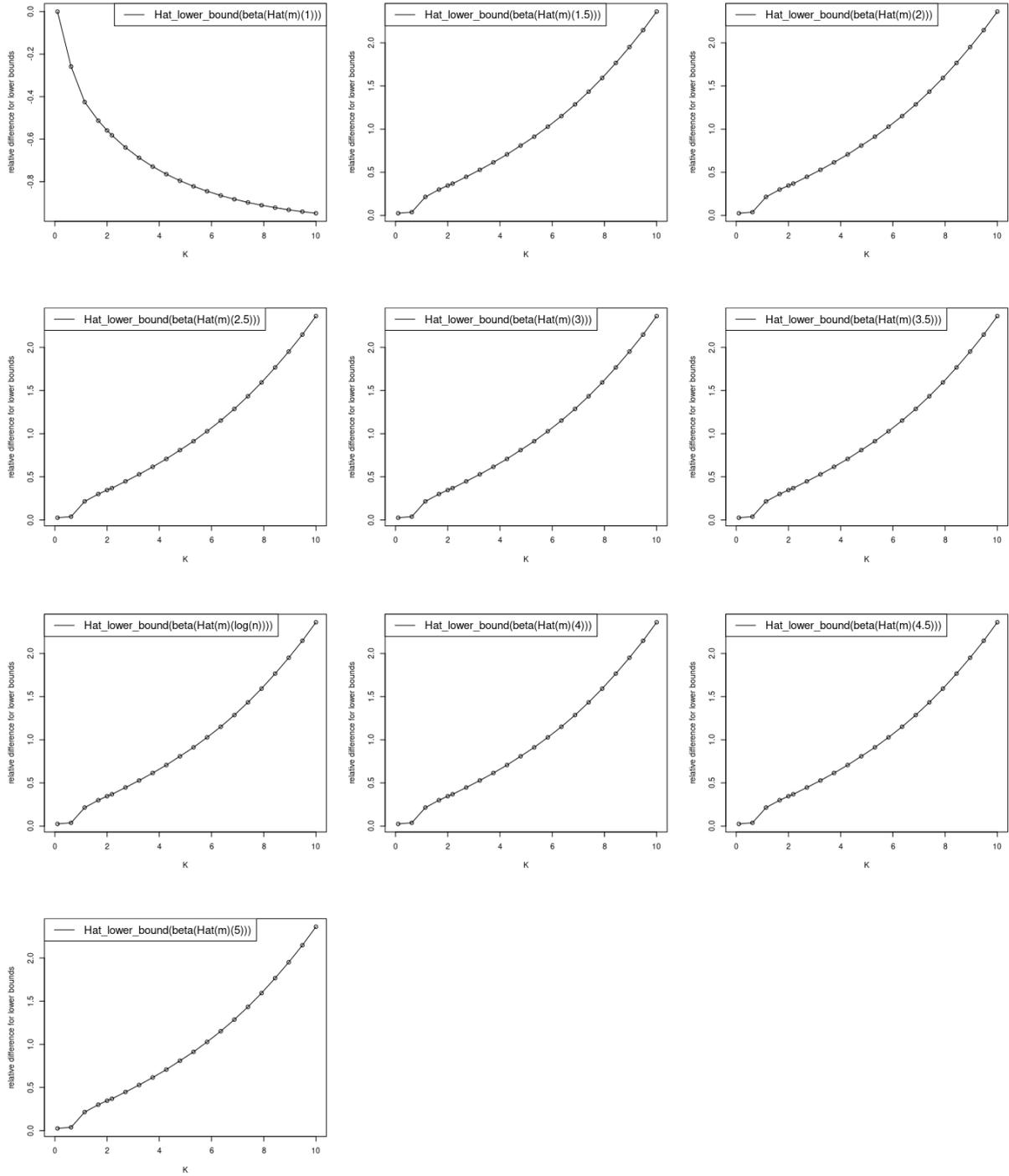


Figure 6: Curves of the relative change values between the function $b(K, \beta^*, \sigma^2)$ and the functions $b(K, \hat{\beta}_{\hat{m}(K)}, \hat{\sigma}^2)$ with respectively $\hat{\beta}_{\hat{m}(1)}$, $\hat{\beta}_{\hat{m}(1.5)}$, $\hat{\beta}_{\hat{m}(2)}$, $\hat{\beta}_{\hat{m}(2.5)}$, $\hat{\beta}_{\hat{m}(3)}$, $\hat{\beta}_{\hat{m}(3.5)}$, $\hat{\beta}_{\hat{m}(4)}$, $\hat{\beta}_{\hat{m}(4.5)}$, $\hat{\beta}_{\hat{m}(5)}$ and $\hat{\beta}_{\hat{m}(\log(n))}$, where estimators are calculating from only one dataset.

Trade-off between prediction and FDR for variable selection

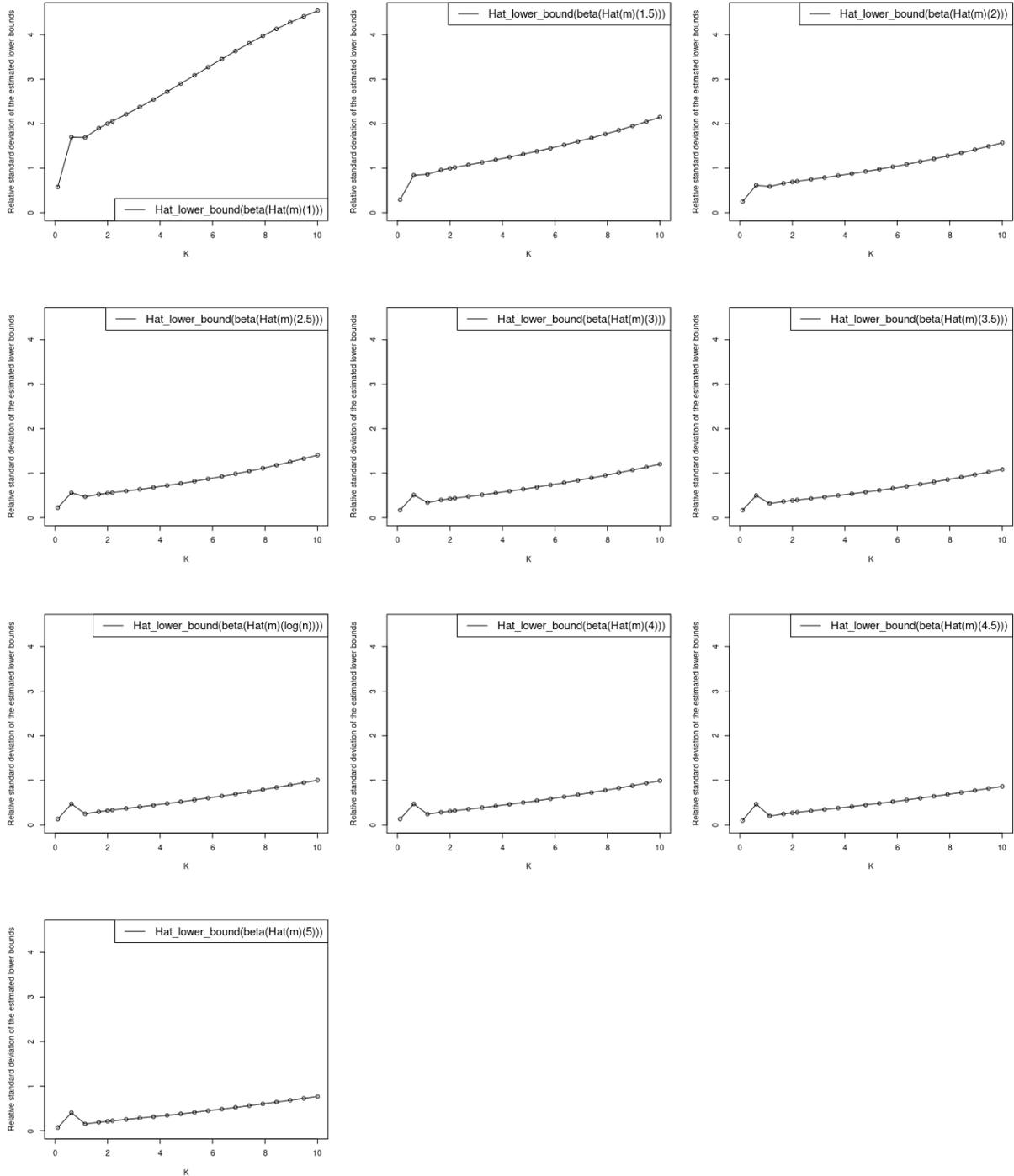


Figure 7: Curves of the relative standard deviation (standard deviation normalized by the mean) of the functions $b(K, \hat{\beta}_{\hat{m}(\bar{K})}, \hat{\sigma}^2)$ obtained from 100 data sets. With each one, $\hat{\beta}_{\hat{m}(1)}, \hat{\beta}_{\hat{m}(1.5)}, \hat{\beta}_{\hat{m}(2)}, \hat{\beta}_{\hat{m}(2.5)}, \hat{\beta}_{\hat{m}(3)}, \hat{\beta}_{\hat{m}(3.5)}, \hat{\beta}_{\hat{m}(4)}, \hat{\beta}_{\hat{m}(4.5)}, \hat{\beta}_{\hat{m}(5)}$ and $\hat{\beta}_{\hat{m}(\log(n))}$ are calculated given $b(K, \hat{\beta}_{\hat{m}(\bar{K})}, \hat{\sigma}^2)$, variance of the 100 $b(K, \hat{\beta}_{\hat{m}(\bar{K})}, \hat{\sigma}^2)$ functions and then the relative standard deviation with respect to K .

Trade-off between prediction and FDR for variable selection

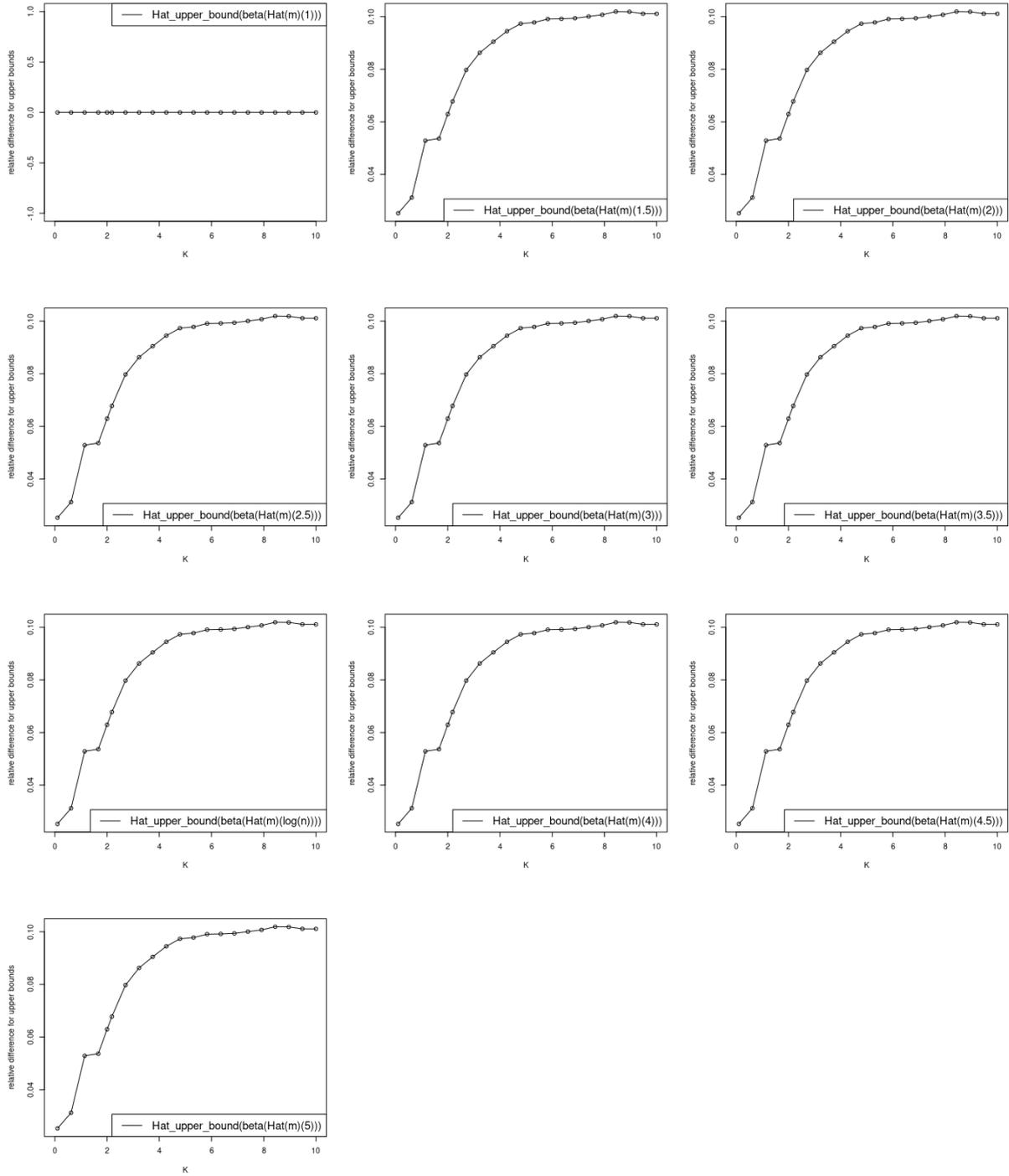


Figure 8: Curves of the relative change values between the function $B(K, \beta^*, \sigma^2)$ and the functions $B(K, \hat{\beta}_{\hat{m}(K)}, \hat{\sigma}^2)$ with respectively $\hat{\beta}_{\hat{m}(1)}$, $\hat{\beta}_{\hat{m}(1.5)}$, $\hat{\beta}_{\hat{m}(2)}$, $\hat{\beta}_{\hat{m}(2.5)}$, $\hat{\beta}_{\hat{m}(3)}$, $\hat{\beta}_{\hat{m}(3.5)}$, $\hat{\beta}_{\hat{m}(4)}$, $\hat{\beta}_{\hat{m}(4.5)}$, $\hat{\beta}_{\hat{m}(5)}$ and $\hat{\beta}_{\hat{m}(\log(n))}$, where estimators are calculating from only one dataset.

Trade-off between prediction and FDR for variable selection

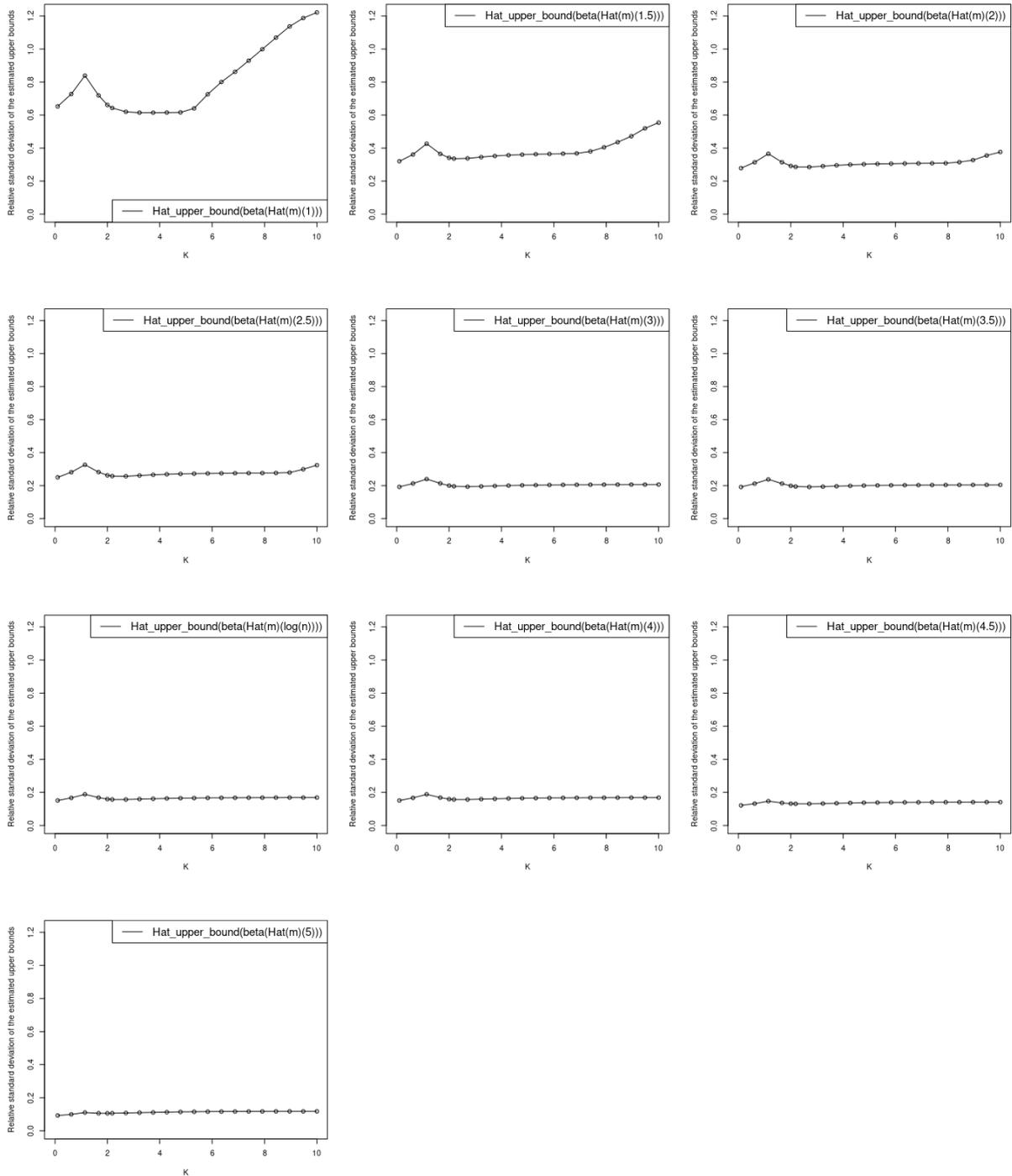


Figure 9: Curves of the relative standard deviation (standard deviation normalized by the mean) of the functions $B(K, \hat{\beta}_{\hat{m}(\bar{K})}, \hat{\sigma}^2)$ obtained from 100 data sets. With each one, $\hat{\beta}_{\hat{m}(1)}, \hat{\beta}_{\hat{m}(1.5)}, \hat{\beta}_{\hat{m}(2)}, \hat{\beta}_{\hat{m}(2.5)}, \hat{\beta}_{\hat{m}(3)}, \hat{\beta}_{\hat{m}(3.5)}, \hat{\beta}_{\hat{m}(4)}, \hat{\beta}_{\hat{m}(4.5)}, \hat{\beta}_{\hat{m}(5)}$ and $\hat{\beta}_{\hat{m}(\log(n))}$ are calculated given $B(K, \hat{\beta}_{\hat{m}(\bar{K})}, \hat{\sigma}^2)$, variance of the 100 $B(K, \hat{\beta}_{\hat{m}(\bar{K})}, \hat{\sigma}^2)$ functions and then the relative standard deviation with respect to K .