

Group SLOPE - adaptive selection of groups of predictors ¹

Damian Brzyski^{a,b}, Alexej Gossman^c, Weijie Su^d, Małgorzata Bogdan^e

^a *Department of Epidemiology and Biostatistics, Indiana University, Bloomington, IN 47405, USA*

^b *Institute of Mathematics, Jagiellonian University, 30-348 Cracow, Poland*

^c *Department of Mathematics, Tulane University, New Orleans, LA 70118, USA*

^d *Department of Statistics, University of Pennsylvania, Philadelphia, PA 19104, USA*

^e *Institute of Mathematics, University of Wrocław, 50-384 Wrocław, Poland*

Key words: Asymptotic Minimax, False Discovery Rate, Group selection, Model Selection, Multiple Regression, SLOPE

Abstract

Sorted L-One Penalized Estimation (SLOPE, [10]) is a relatively new convex optimization procedure which allows for adaptive selection of regressors under sparse high dimensional designs. Here we extend the idea of SLOPE to deal with the situation when one aims at selecting whole groups of explanatory variables instead of single regressors. Such groups can be formed by clustering strongly correlated predictors or groups of dummy variables corresponding to different levels of the same qualitative predictor. We formulate the respective convex optimization problem, gSLOPE (group SLOPE), and propose an efficient algorithm for its solution. We also define a notion of the group false discovery rate (gFDR) and provide a choice of the sequence of tuning parameters for gSLOPE so that gFDR is provably controlled at a prespecified level if the groups of variables are orthogonal to each other. Moreover, we prove that the resulting procedure adapts to unknown sparsity and is asymptotically minimax with respect to the estimation of the proportions of variance of the response variable explained by regressors from different groups. We also provide a method for the choice of the regularizing sequence when variables in different groups are not orthogonal but statistically independent and illustrate its good properties with computer simulations. Finally, we illustrate the advantages of gSLOPE in the context of Genome Wide Association Studies. R package `grpsLOPE` with implementation of our method is available on CRAN.

1 Introduction

Consider the classical multiple regression model of the form

$$y = X\beta + z, \tag{1.1}$$

where y is the n dimensional vector of values of the response variable, X is the n by p experiment (design) matrix and $z \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$. We assume that y and X are known, while β is unknown. In many applications the purpose of the statistical analysis is to recover the support of β , which identifies the set of important regressors. Here, the true support corresponds to truly relevant variables (i.e. variables which have impact on observations). Common procedures to solve this model selection problem rely on minimization of some objective function consisting of the weighted sum of two components: first term responsible for the goodness of fit and second term penalizing the model complexity. Among such procedures one can mention classical model selection criteria like the Akaike Information Criterion (AIC) [3] and the Bayesian Information Criterion (BIC) [22], where the penalty depends on the number of variables included in the model, or LASSO [27], where the penalty depends on the ℓ_1 norm of regression coefficients. The main advantage of LASSO over classical model selection criteria is that it is a convex optimization problem and, as such, it can be easily solved even for very large design matrices.

¹An earlier version of the paper appeared on arXiv.org in November 2015: arXiv:1511.09078

LASSO solution is obtained by solving the optimization problem

$$\arg \min_b \left\{ \frac{1}{2} \|y - Xb\|^2 + \lambda_L \|b\|_1 \right\}, \quad (1.2)$$

where λ_L is a tuning parameter defining the trade-off between the model fit and the sparsity of solution. In practical applications the selection of good λ_L might be very challenging. For example it has been reported that in high dimensional settings the popular cross-validation typically leads to detection of a large number of false regressors (see e.g. [10]). The general rule is that when one reduces λ_L , then LASSO can identify more elements from the true support (true discoveries) but at the same time it generates more false discoveries. In general the numbers of true and false discoveries for a given λ_L depend on unknown properties on the data generating mechanism, like the number of true regressors and the magnitude of their effects. A very similar problem occurs when selecting thresholds for individual tests in the context of multiple testing. Here it was found that the popular Benjamini-Hochberg rule (BH, [7]), aimed at control of the False Discovery Rate (FDR), adapts to the unknown data generating mechanism and has some desirable optimality properties under a variety of statistical settings (see e.g. [1, 8, 20, 15]). The main property of this rule is that it relaxes the thresholds along the sequence of test statistics, sorted in the decreased order of magnitude. Recently the same idea was used in a new generalization of LASSO, named SLOPE (Sorted L-One Penalized Estimation, [9, 10]). Instead of the ℓ_1 norm (as in LASSO case), the method uses FDR control properties of J_λ norm, defined as follows; for sequence $\{\lambda\}_{i=1}^p$ satisfying $\lambda_1 \geq \dots \geq \lambda_p \geq 0$ and $b \in \mathbb{R}^p$, $J_\lambda(b) := \sum_{i=1}^p \lambda_i |b|_{(i)}$, where $|b|_{(1)} \geq \dots \geq |b|_{(p)}$ is the vector of sorted absolute values of coordinates of b . SLOPE is the solution to a convex optimization problem

$$\arg \min_b \left\{ \frac{1}{2} \|y - Xb\|^2 + J_\lambda(b) \right\}, \quad (1.3)$$

which clearly reduces to LASSO for $\lambda_1 = \dots = \lambda_p =: \lambda_L$. Similarly as in classical model selection, the support of solution defines the subset of variables estimated as relevant. In [10] it is shown that when the sequence λ corresponds to the decreasing sequence of threshold for BH then SLOPE controls FDR under orthogonal designs, i.e. when $X^T X = \mathbf{I}_n$. Moreover, in [26] it is proved that SLOPE with this sequence of tuning parameters adapts to unknown sparsity and is asymptotically minimax under orthogonal and random Gaussian designs.

In the sequence of examples presented in [9], [10] and [11] it was shown that SLOPE has very desirable properties in terms of FDR control in case when regressor variables are weakly correlated. While there exist other interesting approaches which allow to control FDR under correlated designs (e.g., [5]), the efforts to prevent detection of false regressors which are strongly correlated with true ones inevitably lead to a loss of power. An alternative approach to deal with strongly correlated predictors is to simply give up the idea of distinguishing between them and include all of them into the selected model as a group. This leads to the problem of group selection in linear regression, extensively investigated and applied in many fields of science. In many of these applications the groups are selected not only due to the strong correlations but also by taking into account the problem specific scientific knowledge. It is also common to cluster dummy variables corresponding to different levels of qualitative predictors.

Probably the most known convex optimization method for selection of groups of explanatory variables is the group LASSO (gLASSO) [4]. For a fixed tuning parameter, $\lambda_{gL} > 0$, the gLASSO estimate is most frequently (e.g. [29], [23]) defined as a solution to optimization problem

$$\arg \min_b \left\{ \frac{1}{2} \left\| y - \sum_{i=1}^m X_{I_i} b_{I_i} \right\|_2^2 + \sigma \lambda_{gL} \sum_{i=1}^m \sqrt{|I_i|} \|b_{I_i}\|_2 \right\}, \quad (1.4)$$

where the sets I_1, \dots, I_m form a partition of the set $\{1, \dots, p\}$, $|I_i|$ denotes the number of elements in set I_i , X_{I_i} is the submatrix of X composed of columns indexed by I_i and b_{I_i} is the restriction

of b to indices from I_i . The method introduced in this article is, however, closer to the alternative version of gLASSO, in which penalties are imposed on $\|X_{I_i}b_{I_i}\|_2$ rather than $\|b_{I_i}\|_2$. This method was formulated in [24], where authors defined estimate of β as

$$\beta^{\text{gL}} := \arg \min_b \left\{ \frac{1}{2} \left\| y - \sum_{i=1}^m X_{I_i} b_{I_i} \right\|_2^2 + \sigma \lambda_{gL} \sum_{i=1}^m \sqrt{|I_i|} \|X_{I_i} b_{I_i}\|_2 \right\}, \quad (1.5)$$

with the condition $\|X_{I_i} \beta_{I_i}^{\text{gL}}\|_2 > 0$ serving as a group relevance indicator.

Similarly as in the context of regular model selection, the properties of gLASSO strongly depend on the shrinkage parameter λ_{gL} , whose optimal value is the function of unknown parameters of true data generating mechanism. Thus, a natural question arises if the idea of SLOPE can be used for construction of the similar adaptive procedure for the group selection. To answer this query in this paper we define and investigate the properties of the group SLOPE (gSLOPE). We formulate the respective optimization problem and provide the algorithm for its solution. We also define the notion of the group FDR (gFDR), and provide the theoretical choice of the sequence of regularization parameters, which guarantees that gSLOPE controls gFDR in the situation when variables in different groups are orthogonal to each other. Moreover, we prove that the resulting procedure adapts to unknown sparsity and is asymptotically minimax with respect to the estimation of the proportions of variance of the response variable explained by regressors from different groups. Additionally, we provide a way of constructing the sequence of regularization parameters under the assumption that the regressors from distinct groups are independent and use computer simulations to show that it allows to control gFDR. Good properties of group SLOPE are illustrated using the practical example of Genome Wide Association Study. R package `grpSLOPE` with implementation of our method is available on CRAN.

2 Group SLOPE

2.1 Formulation of the optimization problem

Let the design matrix X belong to the space $M(n, p)$ of matrices with n rows and p columns. Furthermore, suppose that $I = \{I_1, \dots, I_m\}$ is some partition of the set $\{1, \dots, p\}$, i.e. I_i 's are nonempty sets, $I_i \cap I_j = \emptyset$ for $i \neq j$ and $\bigcup I_i = \{1, \dots, p\}$. We will consider the linear regression model with m groups of the form

$$y = \sum_{i=1}^m X_{I_i} \beta_{I_i} + z, \quad (2.1)$$

where X_{I_i} is the submatrix of X composed of columns indexed by I_i and β_{I_i} is the restriction of β to indices from the set I_i . We will use notations l_1, \dots, l_m to refer to the ranks of submatrices X_{I_1}, \dots, X_{I_m} . To simplify notations in further part, we will assume that $l_i > 0$ (i.e. there is at least one nonzero entry of X_{I_i} for all i). Besides this, X may be absolutely arbitrary matrix, in particular any linear dependencies inside submatrices X_{I_i} are allowed.

In this article we will treat the value $\|X_{I_i} \beta_{I_i}\|_2$ as a measure of an impact of i th group on the response and we will say that the group i is truly relevant if and only if $\|X_{I_i} \beta_{I_i}\|_2 > 0$. Thus our task of the identification of the relevant groups is equivalent with finding the support of the vector $\llbracket \beta \rrbracket_{X, I} := (\|X_{I_1} \beta_{I_1}\|_2, \dots, \|X_{I_m} \beta_{I_m}\|_2)^\top$.

To estimate the nonzero coefficients of $\llbracket \beta \rrbracket_{X, I}$, we will use a new penalized method, namely group SLOPE (gSLOPE). For a given nonincreasing sequence of nonnegative tuning parameters, $\lambda_1, \dots, \lambda_m$, a given sequence of positive weights, w_1, \dots, w_m , and a design matrix, X , the gSLOPE, β^{gS} , is defined as solutions to

$$\beta^{\text{gS}} := \arg \min_b \left\{ \frac{1}{2} \|y - Xb\|_2^2 + \sigma J_\lambda(W \llbracket b \rrbracket_{X, I}) \right\}, \quad (2.2)$$

where W is a diagonal matrix with $W_{i,i} := w_i$, for $i = 1, \dots, m$. The estimate of $\llbracket \beta \rrbracket_{X,I}$ support is simply defined by the indices corresponding to nonzeros of $\llbracket \beta^{\text{gs}} \rrbracket_{X,I}$.

It is easy to see that when one considers p groups containing only one variable (i.e. singleton groups situation), then taking all weights equal to one reduces (2.2) to SLOPE (1.3). On the other hand, taking $w_i = \sqrt{|I_i|}$ and putting $\lambda_1 = \dots = \lambda_m =: \lambda_{gL}$, immediately gives gLASSO problem (1.5) with the smoothing parameter λ_{gL} . The gSLOPE could be therefore treated both: as the extension to SLOPE, and the extension to group LASSO.

Now, let us define $\tilde{p} = l_1 + \dots + l_m$ and consider the following partition, $\mathbb{I} = \{\mathbb{I}_1, \dots, \mathbb{I}_m\}$, of the set $\{1, \dots, \tilde{p}\}$

$$\mathbb{I}_1 := \{1, \dots, l_1\}, \quad \mathbb{I}_2 := \{l_1 + 1, \dots, l_1 + l_2\}, \quad \dots, \quad \mathbb{I}_m := \left\{ \sum_{j=1}^{m-1} l_j + 1, \dots, \sum_{j=1}^m l_j \right\}. \quad (2.3)$$

Observe that each X_{I_i} can be represented as $X_{I_i} = U_i R_i$, where U_i is a matrix with l_i orthogonal columns of a unit l_2 norm, whose span coincides with the space spanned by the columns of X_{I_i} , and R_i is the corresponding matrix of a full row rank. Define n by l matrix \tilde{X} by putting $\tilde{X}_{\mathbb{I}_i} := U_i$ for $i = 1, \dots, m$. Now observe that denoting $c_{\mathbb{I}_i} := R_i b_{I_i}$ for $i \in \{1, \dots, m\}$ we immediately obtain

$$\begin{aligned} Xb &= \sum_{i=1}^m X_{I_i} b_{I_i} = \sum_{i=1}^m U_i R_i b_{I_i} = \sum_{i=1}^m \tilde{X}_{\mathbb{I}_i} c_{\mathbb{I}_i} = \tilde{X}c, \\ \left(\llbracket b \rrbracket_{X,I} \right)_i &= \|X_{I_i} b_{I_i}\|_2 = \|R_i b_{I_i}\|_2 = \|c_{\mathbb{I}_i}\|_2 \end{aligned} \quad (2.4)$$

and the problem (2.2) can be equivalently presented in the form

$$\begin{cases} c^{\text{gs}} := \arg \min_c \left\{ \frac{1}{2} \|y - \tilde{X}c\|_2^2 + \sigma J_\lambda(W \llbracket c \rrbracket_{\mathbb{I}}) \right\}, \\ c_{\mathbb{I}_i}^{\text{gs}} := R_i \beta_{I_i}^{\text{gs}}, \quad i = 1, \dots, m \end{cases}, \quad (2.5)$$

for $\llbracket c \rrbracket_{\mathbb{I}} := (\|c_{\mathbb{I}_1}\|_2, \dots, \|c_{\mathbb{I}_m}\|_2)^\top$. Therefore to identify the relevant groups and estimate their group effects it is enough to solve the optimization problem (2.5). We will say that (2.5) is the standardized version of the problem (2.2).

Remark 2.1. *Similar formulation of the group SLOPE was proposed in [16]. However [16] considers only the case when the weights w_i are equal to the square root of the group size and penalties are imposed directly on $\|\beta_{I_i}\|_2$ rather than on group effects $\|X_{I_i} \beta_{I_i}\|_2$. This makes the method of [16] dependent on scaling or rotations of variables in a given group. In comparison to [16], where a Monte Carlo approach for estimating the regularizing sequence was proposed, our article provides choice of the smoothing parameters which provably allow for FDR control in case where the regressors in different groups are orthogonal to each other and its modification, which according to our simulation study allows for FDR control where regressors in different groups are independent.*

2.2 Numerical algorithm

As shown in Appendix B the function $J_{\lambda, W, \mathbb{I}}(b) := J_\lambda(W \llbracket b \rrbracket_{\mathbb{I}})$ is a norm and the optimization problem (2.5) can be solved by using proximal gradient methods. In our R package `grpSLOPE` available on CRAN (The Comprehensive R Archive Network) the accelerated proximal gradient method known as FISTA [6] is applied, which uses the specific procedure for choosing steps sizes, to achieve fast convergence rate. The proximal operator for gSLOPE is obtained by appropriate transformation and reduction of the problem, so the fast proximal operator for SLOPE [9] can be used. To derive proper stopping criteria, we have considered dual problem to gSLOPE and employed the strong duality property. The detailed description of the proximal operator for gSLOPE as well as of the dual norm and conjugate of grouped sorted l_1 norm is provided in the Appendix B.

2.3 Group FDR

Group SLOPE is designed to select groups of variables, which might be very strongly correlated within a group or even linearly dependent. In this context we do not intend to identify single important predictors but rather want to point at the groups which contain at least one true regressor. To theoretically investigate the properties of gSLOPE in this context we now introduce the respective notion of group FDR (gFDR).

Definition 2.2. Consider model (2.1) and let β^{gs} be an estimate given by (2.2). We define two random variables: the number of all groups selected by gSLOPE (Rg) and the number of groups falsely discovered by gSLOPE (Vg), as

$$Rg := |\{i : \|X_{I_i}\beta_{I_i}^{\text{gs}}\|_2 \neq 0\}|, \quad Vg := |\{i : \|X_{I_i}\beta_{I_i}\|_2 = 0, \|X_{I_i}\beta_{I_i}^{\text{gs}}\|_2 \neq 0\}|.$$

Definition 2.3. We define the false discovery rate for groups (gFDR) as

$$gFDR := \mathbb{E} \left[\frac{Vg}{\max\{Rg, 1\}} \right]. \quad (2.6)$$

2.4 Control of gFDR when variables from different groups are orthogonal

Our goal is the identification of the regularizing sequence for gSLOPE such that gFDR can be controlled at any given level $q \in (0, 1)$. In this section we will provide such a sequence, which provably controls gFDR in case when variables in different groups are orthogonal to each other. In subsequent sections we will replace this condition with the weaker assumption of the stochastic independence of regressors in different groups. Before the statement of the main theorem on gFDR control, we will recall the definition of χ distribution and define a scaled χ distribution.

Definition 2.4. We will say that a random variable X_1 has a χ distribution with l degrees of freedom, and write $X_1 \sim \chi_l$, when X_1 could be expressed as $X_1 = \sqrt{X_2}$, for X_2 having a χ^2 distribution with l degrees of freedom. We will say that a random variable X_1 has a scaled χ distribution with l degrees of freedom and scale \mathcal{S} , when X_1 could be expressed as $X_1 = \mathcal{S} \cdot X_2$, for X_2 having a χ distribution with l degrees of freedom. We will use the notation $X_1 \sim \mathcal{S}\chi_l$.

Theorem 2.5 (gFDR control under orthogonal case). Consider model (2.1) with the design matrix X satisfying $X_{I_i}^\top X_{I_j} = 0$, for any $i \neq j$. Denote the number of zero coefficients in $\llbracket \beta \rrbracket_{X,I}$ by m_0 and let w_1, \dots, w_m be positive numbers. Moreover, define the sequence of regularizing parameters $\lambda^{\text{max}} = (\lambda_1^{\text{max}}, \dots, \lambda_m^{\text{max}})^\top$, with

$$\lambda_i^{\text{max}} := \max_{j=1, \dots, m} \left\{ \frac{1}{w_j} F_{\chi_{l_j}}^{-1} \left(1 - \frac{q \cdot i}{m} \right) \right\}, \quad (2.7)$$

where $F_{\chi_{l_j}}$ is a cumulative distribution function of χ distribution with l_j degrees of freedom. Then any solution, β^{gs} , to problem gSLOPE (2.2) generates the same vector $\llbracket \beta^{\text{gs}} \rrbracket_{X,I}$ and it holds

$$gFDR = \mathbb{E} \left[\frac{Vg}{\max\{Rg, 1\}} \right] \leq q \cdot \frac{m_0}{m}.$$

Proof. We will start with the standardized version of the gSLOPE problem, given by (2.5). Based on results discussed in Appendix C, we can consider an equivalent formulation of (2.5)

$$\begin{cases} c^* = \arg \min_c \left\{ \frac{1}{2} \sum_{i=1}^m (\|\tilde{y}_{\mathbb{I}_i}\|_2 - w_i^{-1} c_i)^2 + J_{\sigma\lambda}(c) \right\} \\ \|X_{I_i}\beta_{I_i}^{\text{gs}}\|_2 = c_i^* (w_i \|\tilde{y}_{\mathbb{I}_i}\|_2)^{-1} \tilde{y}_{\mathbb{I}_i}, \quad i = 1, \dots, m, \end{cases} \quad (2.8)$$

where $\tilde{y} = \tilde{X}^T y$ has a multivariate normal distribution $\mathcal{N}(\tilde{\beta}, \sigma^2 \mathbf{I}_{\tilde{p}})$ with $\tilde{\beta}_{I_i} = R_i \beta_{I_i}$. The uniqueness of $\llbracket \beta^{\text{gs}} \rrbracket_{X,I}$ follows simply from the uniqueness of c^* in (2.8). Define random variables $R := |\{i : c_i^* \neq 0\}|$ and $V := |\{i : \|\tilde{\beta}_{I_i}\|_2 = 0, c_i^* \neq 0\}|$. Clearly, then $Rg = R$ and $Vg = V$. Consequently, it is enough to show that

$$\mathbb{E} \left[\frac{V}{\max\{R, 1\}} \right] \leq q \cdot \frac{m_0}{m}.$$

Without loss of generality we can assume that groups I_1, \dots, I_{m_0} are truly irrelevant, which gives $\|\tilde{\beta}_{I_1}\|_2 = \dots = \|\tilde{\beta}_{I_{m_0}}\|_2 = 0$ and $\|\tilde{\beta}_{I_j}\|_2 > 0$ for $j > m_0$. Suppose now that r, i are some fixed indices from $\{1, \dots, m\}$. From definition of λ_r^{max}

$$\lambda_r^{\text{max}} \geq \frac{1}{w_i} F_{\chi_{l_i}}^{-1} \left(1 - \frac{qr}{m} \right) \implies 1 - F_{\chi_{l_i}}(\lambda_r^{\text{max}} w_i) \leq \frac{qr}{m}. \quad (2.9)$$

Now, let us fix $i \leq m_0$. Since $\sigma^{-1} \|\tilde{y}_{I_i}\|_2 \sim \chi_{l_i}$ we have

$$\mathbb{P}(w_i^{-1} \|\tilde{y}_{I_i}\|_2 \geq \sigma \lambda_r^{\text{max}}) = \mathbb{P}(\sigma^{-1} \|\tilde{y}_{I_i}\|_2 \geq \lambda_r^{\text{max}} w_i) = 1 - F_{\chi_{l_i}}(\lambda_r^{\text{max}} w_i) \leq \frac{qr}{m}. \quad (2.10)$$

Now, denote by \tilde{R}^i the number of nonzero coefficients in SLOPE estimate (2.8) after eliminating i th group of explanatory variables. Thanks to lemmas D.6 and D.7, we immediately get

$$\{\llbracket \tilde{y} \rrbracket_{\mathbb{I}} : c_i^* \neq 0 \text{ and } R = r\} \subset \{\llbracket \tilde{y} \rrbracket_{\mathbb{I}} : w_i^{-1} \|\tilde{y}_{I_i}\|_2 > \sigma \lambda_r^{\text{max}} \text{ and } \tilde{R}^i = r - 1\}, \quad (2.11)$$

which together with (2.10) raises

$$\begin{aligned} \mathbb{P}(c_i^* \neq 0 \text{ and } R = r) &\leq \mathbb{P}(w_i^{-1} \|\tilde{y}_{I_i}\|_2 > \sigma \lambda_r^{\text{max}} \text{ and } \tilde{R}^i = r - 1) \\ &= \mathbb{P}(w_i^{-1} \|\tilde{y}_{I_i}\|_2 > \sigma \lambda_r^{\text{max}}) \mathbb{P}(\tilde{R}^i = r - 1) \\ &\leq \frac{qr}{m} \mathbb{P}(\tilde{R}^i = r - 1). \end{aligned} \quad (2.12)$$

Therefore

$$\begin{aligned} \mathbb{E} \left[\frac{V}{\max\{R, 1\}} \right] &= \sum_{r=1}^m \mathbb{E} \left[\frac{V}{r} \mathbf{1}_{\{R=r\}} \right] = \sum_{r=1}^m \frac{1}{r} \mathbb{E} \left[\sum_{i=1}^{m_0} \mathbf{1}_{\{c_i^* \neq 0\}} \mathbf{1}_{\{R=r\}} \right] = \\ &\sum_{r=1}^m \frac{1}{r} \sum_{i=1}^{m_0} \mathbb{P}(c_i^* \neq 0 \text{ and } R = r) \leq \sum_{i=1}^{m_0} \frac{q}{m} \sum_{r=1}^m \mathbb{P}(\tilde{R}^i = r - 1) = \frac{qm_0}{m}, \end{aligned} \quad (2.13)$$

which finishes the proof. ■

Figure 1 illustrates the performance of gSLOPE under the design matrix $X = \mathbf{I}_p$ (hence the rank of i group, l_i , coincides with its size), with $p = 5000$. In Figure 1 (a) all groups are of the same size $l = 5$, while in Figures 1 (b) - (d) the explanatory variables are clustered into $m = 1000$ groups of sizes from the set $\{3, 4, 5, 6, 7\}$; 200 groups of each size. Each coefficient of β_{I_i} , in truly relevant group i , was generated independently from $U[0.1, 1.1]$ distribution and then β_{I_i} was scaled such that $(\llbracket \beta \rrbracket_{X,I})_i = a\sqrt{l_i}$. Parameter a was selected to satisfy the condition $\frac{1}{m} \sum_{i=1}^m a\sqrt{l_i} = \frac{1}{m} \sum_{i=1}^m B(m, l_i)$, where $B(m, l)$ is defined in (F.4). Such signals are comparable to the maximal noise and can be detected with moderate power, which allows for a meaningful comparison between different methods.

Figure 1 (a) illustrates that the sequence λ^{max} allows to keep gFDR very close to the "nominal" level when groups are of the same size. However, Figure 1 (b) shows that for groups of different size λ^{max} is rather conservative, i.e. the achieved gFDR is significantly lower than assumed. This

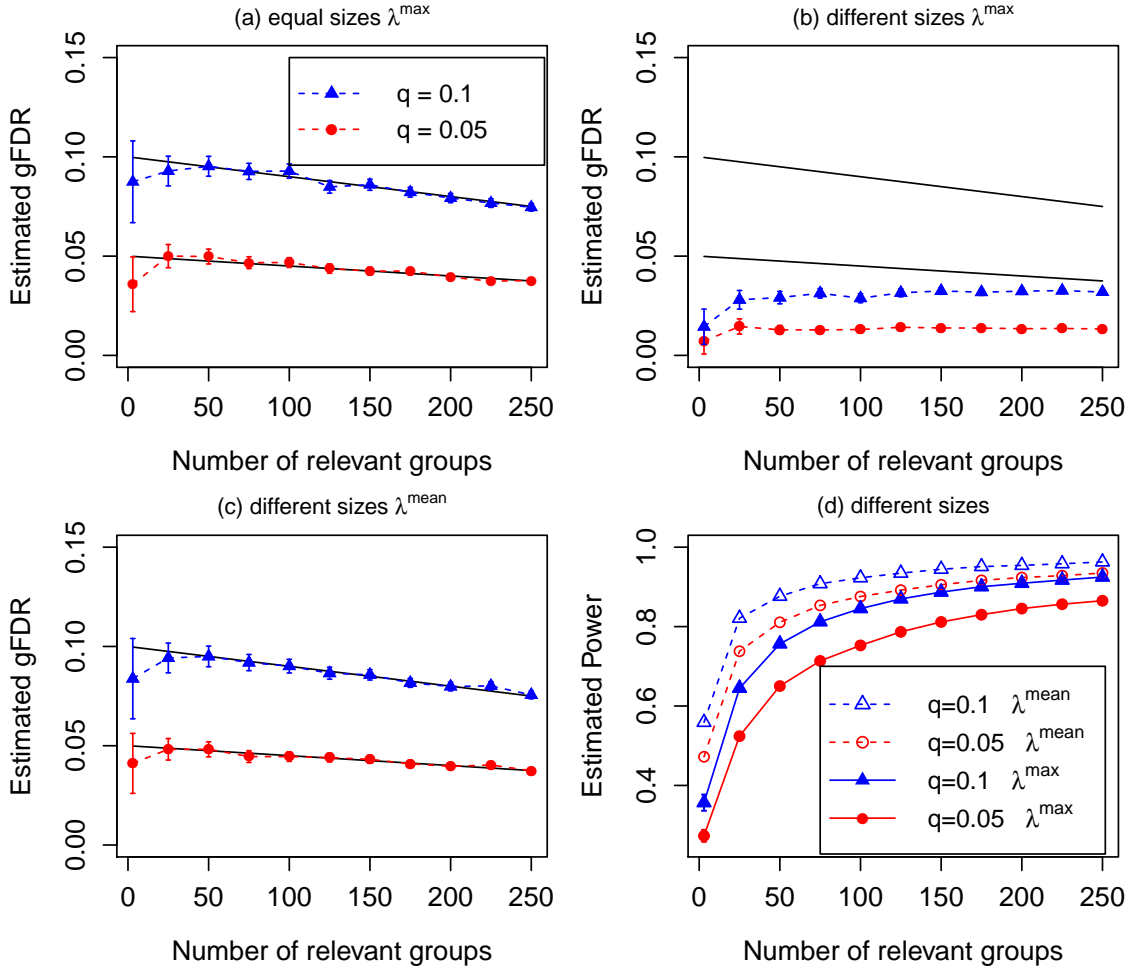


Figure 1: Orthogonal situation with $n = p = 5000$ and $m = 1000$. In (a) all groups are of the same size $l = 5$, while in (b)-(d) there are 200 groups of each of sizes $l_i \in \{3, 4, 5, 6, 7\}$. In (a) and (b) gSLOPE works with the regularizing sequence λ^{\max} , while in (c) and (d) λ^{mean} is used. For each target gFDR level and true support size, 300 iterations were performed. Bars correspond to $\pm 2\text{SE}$. Black straight lines represent the "nominal" gFDR level $q \cdot ((m - k)/m)$, for k being true support size. Weights are defined as $w_i := \sqrt{l_i}$.

suggests that the shrinkage (dictated by λ) could be slightly decreased, such that the method gets more power and still achieves the gFDR below the assumed level. Returning to the proof of Theorem 2.5, we can see that for each $i \in \{1, \dots, m\}$ we have

$$1 - F_{\chi_{l_i}}(\lambda_r^{\max} w_i) \leq \frac{qr}{m}, \quad (2.14)$$

with equality holding only for i being the index of the maximum in (2.7). In the result the inequality in (2.13) is usually strict and the true gFDR might be substantially smaller than the nominal level. The natural relaxation of (2.14) is to require only that

$$\sum_{i=1}^m \left(1 - F_{w_i^{-1}\chi_{l_i}}(\lambda_r)\right) \leq qr. \quad (2.15)$$

Replacing the inequality in (2.15) by equality yields the strategy of choosing the relaxed λ sequence

$$\lambda_r^{mean} := \bar{F}^{-1} \left(1 - \frac{qr}{m} \right) \quad \text{for} \quad \bar{F}(x) := \frac{1}{m} \sum_{i=1}^m F_{w_i^{-1}\chi_{l_i}}(x), \quad r \in \{1, \dots, m\}, \quad (2.16)$$

where $F_{w_i^{-1}\chi_{l_i}}$ is the cumulative distribution function of scaled chi distribution with l_i degrees of freedom and scale $\mathcal{S} = w_i^{-1}$. In Figure 1c we present estimated gFDR, for tuning parameters given by (2.16). The results suggest that with relaxed version of tuning parameters, we can still achieve the "average" gFDR control, where the "average" is with respect to the uniform distribution over all possible signal placements. As shown in Figure 1d, application of λ^{mean} allows to achieve a substantially larger power than the one provided by λ^{max} . Such a strategy could be especially important in situation, when differences between the smallest and the largest quantiles (among distributions $w_i^{-1}\chi_{l_i}$) are relatively large and all groups have the same prior probability of being relevant.

2.5 The accuracy of estimation

Up until this point, we have only considered the testing properties of gSLOPE. Though originally proposed to control the FDR, surprisingly, SLOPE enjoys appealing estimation properties as well [26]. It thus would be desirable to extend this link between testing and estimation for gSLOPE. In measuring the deviation of an estimator from the ground truth β , as earlier, we focus on the group level instead of an individual. Accordingly, here we aim to estimate parts of variance of Y explained by every group, which are contained in the vector $[\beta]_{X,I} := (\|X_{I_1}\beta_{I_1}\|_2, \dots, \|X_{I_m}\beta_{I_m}\|_2)^\top$ or $[\tilde{\beta}]_{\mathbb{I}} := (\|\tilde{\beta}_{\mathbb{I}_1}\|_2, \dots, \|\tilde{\beta}_{\mathbb{I}_m}\|_2)^\top$, equivalently. For illustration purpose, we employ the setting described as follows. Imagine that we have a sequence of problems with the number of groups m growing to infinity: the design X is orthonormal at groups level; ranks of submatrices X_{I_i} , l_i , are bounded, that is, $\max l_i \leq l$ for some constant integer l ; denoting by $k \geq 1$ the sparsity level (that is, the number of relevant groups), we assume the asymptotics $k/m \rightarrow 0$. Now we state our minimax theorem, where we write $a \sim b$ if $a/b \rightarrow 1$ in the asymptotic limit, and $\|[\beta]_{X,I}\|_0$ denotes the number of nonzero entries of $[\beta]_{X,I}$. The proof makes use of the same techniques for proving Theorem 1.1 in [26] and is deferred to the Appendix.

Theorem 2.6. *Fix any constant $q \in (0, 1)$, let $w_i = 1$ and $\lambda_i = F_{\chi_{l_i}}^{-1}(1 - qi/m)$ for $i = 1, \dots, m$. Under the preceding conditions, gSLOPE is asymptotically minimax over the nearly black object $\{\beta : \|[\beta]_{X,I}\|_0 \leq k\}$, i.e.,*

$$\sup_{\|[\beta]_{X,I}\|_0 \leq k} \mathbb{E} \left(\left\| [\beta^{gs}]_{X,I} - [\beta]_{X,I} \right\|_2^2 \right) \sim \inf_{\hat{\beta}} \sup_{\|[\beta]_{X,I}\|_0 \leq k} \mathbb{E} \left(\left\| [\hat{\beta}]_{X,I} - [\beta]_{X,I} \right\|_2^2 \right),$$

where the infimum is taken over all measurable estimators $\hat{\beta}(y, X)$.

Notably, in this theorem the choice of λ_i does not assume the knowledge of sparsity level. Or putting it differently, in stark contrast to gLASSO, gSLOPE is adaptive to a range of sparsity in achieving the exact minimaxity. Combining Theorems 2.5 and 2.6, we see the remarkable link between FDR control and minimax estimation also applies to gSLOPE [1, 26]. While it is out of the scope of this paper, it is of great interest to extend this minimax result to general design matrices.

2.6 The impact of chosen weights

In this subsection we will discuss the influence of chosen weights, $\{w_i\}_{i=1}^m$, on results. Let $I = \{I_1, \dots, I_m\}$ be a given partition into groups and l_1, \dots, l_m be ranks of submatrices X_{I_i} . Assume

the orthogonality at group level, i.e., that it holds $X_{I_i}^\top X_{I_j} = 0$, for $i \neq j$, and suppose that $\sigma = 1$. The support of $\llbracket \beta \rrbracket_{X,I}$ coincides with the support of vector c^* defined in (2.8), namely

$$c^* = \arg \min_c \frac{1}{2} \left\| \llbracket \tilde{y} \rrbracket_{\mathbb{I}} - W^{-1}c \right\|_2^2 + J_\lambda(c), \quad (2.17)$$

where W^{-1} is a diagonal matrix with positive numbers $w_1^{-1}, \dots, w_m^{-1}$ on the diagonal. Suppose now, that c^* has exactly r nonzero coefficients. From Corollary D.4, these indices are given by $\{\pi(1), \dots, \pi(r)\}$, where π is permutation which orders $W^{-1} \llbracket \tilde{y} \rrbracket_{\mathbb{I}}$. Hence, the order of realizations $\{w_i^{-1} \|\tilde{y}_{\mathbb{I}_i}\|_2\}_{i=1}^m$ decides about the subset of groups labeled by gSLOPE as relevant. Suppose that groups I_i and I_j are truly relevant, i.e., $\|\tilde{\beta}_{\mathbb{I}_i}\|_2 > 0$ and $\|\tilde{\beta}_{\mathbb{I}_j}\|_2 > 0$. The distributions of $\|\tilde{y}_{\mathbb{I}_i}\|_2$ and $\|\tilde{y}_{\mathbb{I}_j}\|_2$ are noncentral χ distributions, with l_i and l_j degrees of freedom, and the noncentrality parameters equal to $\|\tilde{\beta}_{\mathbb{I}_i}\|_2$ and $\|\tilde{\beta}_{\mathbb{I}_j}\|_2$, respectively. Now, the expected value of the noncentral χ distribution could be well approximated by the square root of the expected value of the noncentral χ^2 distribution, which gives

$$\mathbb{E}(w_i^{-1} \|\tilde{y}_{\mathbb{I}_i}\|_2) \approx w_i^{-1} \sqrt{\mathbb{E}(\|\tilde{y}_{\mathbb{I}_i}\|_2^2)} = w_i^{-1} \sqrt{l_i + \|\tilde{\beta}_{\mathbb{I}_i}\|_2^2}.$$

Therefore, roughly speaking, truly relevant groups I_i and I_j are treated as comparable, when it occurs $l_i/w_i^2 + \|\tilde{\beta}_{\mathbb{I}_i}\|_2^2/w_i^2 \approx l_j/w_j^2 + \|\tilde{\beta}_{\mathbb{I}_j}\|_2^2/w_j^2$. This gives us the intuition about the behavior of gSLOPE with the choice $w_i = \sqrt{l_i}$ for each i . Firstly, gSLOPE treats all irrelevant groups as comparable, i.e. the size of the group has a relatively small influence on it being selected as a false discovery. Secondly, gSLOPE treats two truly relevant groups as comparable, if groups effect sizes satisfy the condition $(\llbracket \beta \rrbracket_{X,I})_i / (\llbracket \beta \rrbracket_{X,I})_j \approx \sqrt{l_i} / \sqrt{l_j}$. The derived condition could be recast as $\|X_{I_i} \beta_{I_i}\|_2^2 / l_i \approx \|X_{I_j} \beta_{I_j}\|_2^2 / l_j$. This gives a nice interpretation: with the choice $w_i := \sqrt{l_i}$, gSLOPE treats two groups as comparable, when these groups have similar squared effect group sizes per coefficient. One possible idealistic situation, when such a property occurs, is when all β_i 's in truly relevant groups are comparable.

In Figure 2 we see that when the condition $(\llbracket \beta \rrbracket_{X,I})_i / (\llbracket \beta \rrbracket_{X,I})_j = \sqrt{l_i} / \sqrt{l_j}$ is met, the fractions of groups with different sizes in the selected truly relevant groups (STRG) are approximately equal. To investigate the impact of selected weights on the set of discovered groups, we performed simulations with different settings, namely we used $w_i = 1$ and $w_i = l_i$ (without changing other parameters). With the first choice, larger groups are penalized less than before, while the second choice yields the opposite situation. This is reflected in the proportion of each groups in STRG (Figure 2). The values of gFDR are very similar under all choices of weights.

2.7 Independent groups and unknown σ

The assumption that variables in different groups are orthogonal to each other can be satisfied only in rare situations of specifically designed experiments. However, in a variety of applications one can assume that variables in different groups are independent. Such a situation occurs for example in the context of identifying influential genes using distant genetic markers, whose genotypes can be considered as stochastically independent. In this case a group can be formed by clustering dummy variables corresponding to different genotypes of a given marker. Though the difference between stochastic independence and algebraic orthogonality seems rather small, it turns out that small sample correlations between independent regressors together with the shrinkage of regression coefficients lead to magnifying the effective noise and require the adjustment of the tuning sequence λ (see [25] for discussion of this phenomenon in the context of LASSO). Concerning regular SLOPE, this problem was addressed by heuristic modification of λ , proposed in [10] and [9]. This modified sequence was calculated based upon the assumption that explanatory variables are randomly sampled from the Gaussian distribution. However, simulation results from [9] illustrate that it

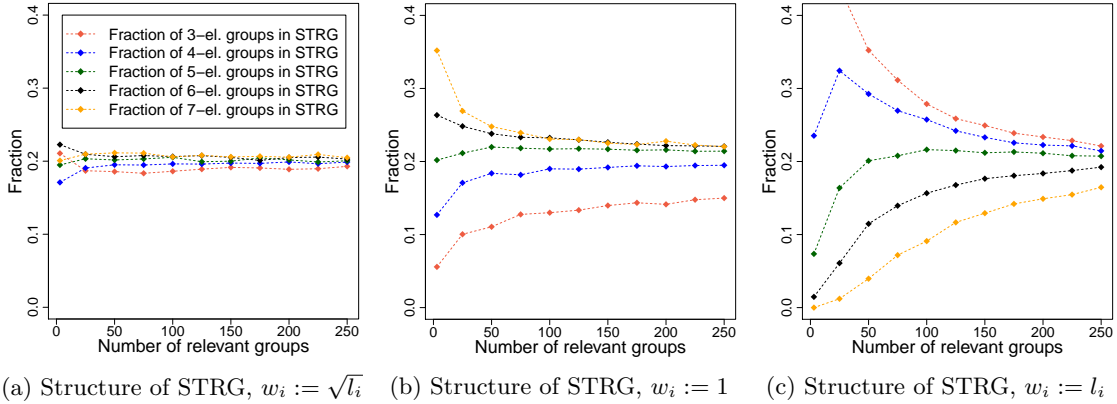


Figure 2: Fraction of each group sizes in selected truly relevant groups (STRG). Beyond the weights, this simulation was conducted with the same setting as in experiments summarized in Figure 1 for λ^{mean} . In particular, for truly relevant groups i and j , it occurs $(\|\beta\|_{X,I})_i / (\|\beta\|_{X,I})_j = \sqrt{l_i} / \sqrt{l_j}$. Target gFDR level was fixed as 0.05.

controls FDR also in case when the columns of the design matrix correspond to additive effects of independent SNPs and the number of causal genes is moderately small.

Following ideas for regular SLOPE presented in [9], we propose the Procedure 1 for calculating the sequence of tuning parameters in case when variables in different groups are independent. The heuristics justifying this choice are substantially more technically involved than the heuristics for regular SLOPE and their details are presented in the Appendix G. Procedure 1 is based on the sequence of λ^{mean} but the version for the conservative choice λ^{max} follows analogously. The proposed sequence of tuning parameters flattens out for a certain value k^* dependent on q, n and l_1, \dots, l_m . It is supposed to control gFDR when the number of identified groups is not much larger than k^* .

Procedure 1 Sequence of tuning parameters for independent groups

input: $q \in (0, 1)$, $w_1, \dots, w_m > 0$, $p, n, m, l_1, \dots, l_m \in \mathbb{N}$
 $\lambda_i := \bar{F}^{-1}(1 - \frac{q}{m})$, for $\bar{F}(x) := \frac{1}{m} \sum_{i=1}^m F_{w_i^{-1} \chi_{l_i}}(x)$;
for $i \in \{2, \dots, m\}$:
 $\lambda^S := (\lambda_1, \dots, \lambda_{i-1})^\top$;
 $\mathcal{S}_j := \sqrt{\frac{n - l_j(i-1)}{n} + \frac{w_j^2 \|\lambda^S\|_2^2}{n - l_j(i-1) - 1}}$, for $j \in \{1, \dots, m\}$;
 $\lambda_i^* := \bar{F}_S^{-1}(1 - \frac{qi}{m})$, for $\bar{F}_S(x) := \frac{1}{m} \sum_{j=1}^m F_{\mathcal{S}_j w_j^{-1} \chi_{l_j}}(x)$;
if $\lambda_i^* \leq \lambda_{i-1}$, then put $\lambda_i := \lambda_i^*$. Otherwise, stop the procedure and put $\lambda_j := \lambda_{i-1}$ for $j \geq i$;
end for

Up until this moment, we have used σ in gSLOPE optimization problem, assuming that this parameter is known. However, in many applications σ is unknown and its estimation is an important issue. When $n > p$, the standard procedure is to use the unbiased estimator of σ^2 , $\hat{\sigma}_{OLS}^2$, given by

$$\hat{\sigma}_{OLS}^2 := (y - X\beta^{OLS})^\top (y - X\beta^{OLS}) / (n - p), \text{ for } \beta^{OLS} := (X^\top X)^{-1} X^\top y. \quad (2.18)$$

For the target situation, with p much larger than n , such an estimator can not be used. To estimate σ we will therefore apply the procedure which was dedicated for this purpose in [9] in the context of SLOPE. Below we present algorithm adjusted to gSLOPE (Procedure 2). The idea standing behind the procedure is simple. The gSLOPE property of producing sparse estimators is used, and in each

Procedure 2 gSLOPE with estimation of σ

input: y , X and λ (defined for some fixed q)
initialize: $S_+ = \emptyset$;
repeat
 $S = S_+$;
 compute RSS obtained by regressing y onto variables in S ;
 set $\hat{\sigma}^2 = RSS/(n - |S| - 1)$;
 compute the solution β^{gS} to gSLOPE with parameters $\hat{\sigma}$ and sequence λ ;
 set $S_+ = \text{supp}(\beta^{gS})$;
until $S_+ = S$

iteration columns in design matrix are first restricted to support of β^{gS} , so that the number of rows exceeds the number of columns and (2.18) can be used. Algorithm terminates when gSLOPE finds the same subset of relevant variables as in the preceding iteration.

To investigate the performance of gSLOPE under the Gaussian design and various group sizes, we performed simulations with 1000 groups. Their sizes were drawn from the binomial distribution, $Bin(1000; 0.008)$, so as the expected value of the group size was equal to 8 (Figure 3c). As a result, we obtained 7917 variables, divided into 1000 groups (the same division was used in all iterations and scenarios). For each sparsity level and the gFDR level 0.1, and each iteration we generated entries of the design matrix using $\mathcal{N}(0, \frac{1}{n})$ distribution, then X was standardized and the values of response variable were generated according to model (2.1) with $\sigma = 1$ and signals generated as in simulations for Figure 1. To identify relevant groups based on the simulated data we have used the iterative version of gSLOPE, with σ estimation (Procedure 2) and lambdas given by Procedure 1. We performed 200 repetitions for each scenario, n was fixed as 5000. Results are represented in Figure 3 and show that our procedure allows to control gFDR at the assumed level.

Additionally, Figure 3 compares gSLOPE to gLASSO with two choices of the smoothing parameter λ . Firstly, we used $\lambda = \lambda_1^{mean}$, which allows to control FDR under the total null hypothesis. Secondly, for each of the iterations we chose λ based on leave-one-out cross-validation. It turns out that the first of these choices becomes rather conservative when the number of truly relevant groups increases. Then gLASSO has a smaller FDR but also a much smaller power than SLOPE (by a factor of three for $k = 60$). Cross-validation works in the opposite way - it yields a large power but also results in a huge proportion of false discoveries, which in our simulations systematically exceeds 60%.

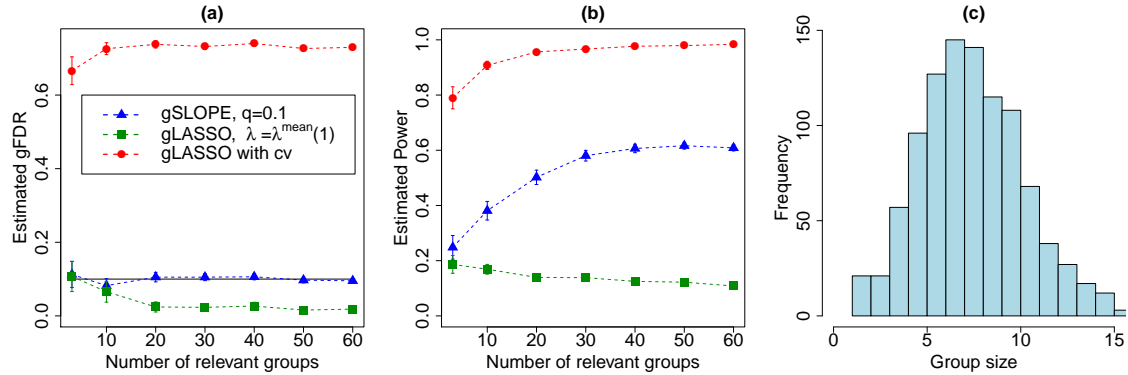


Figure 3: Independent regressors and various group sizes: $m = 1000$, $p = 7917$ and $n = 5000$. Bars correspond to $\pm 2SE$. Entries of design matrix were drawn from $\mathcal{N}(0, 1/n)$ distribution and truly relevant signal, i , was generated such as $\|X_{I_i} \beta_{I_i}\|_2 = \frac{1}{m} \sum_{i=1}^m B(m, l_i)$, where $B(m, l)$ is defined in (F.4).

Table 1: Coding for explanatory variables

genotype	additive dummy variable \tilde{X}	dominance dummy variable \tilde{Z}
<i>aa</i>	2	0
<i>aA</i>	1	1
<i>AA</i>	0	0

2.8 Simulations in the context of Genome-Wide Association Studies

To test the performance of gSLOPE in the context of Genome-Wide Association Studies (GWAS) we have used the North Finland Birth Cohort (NFBC) dataset, available in dbGaP with accession number phs000276.v2.p1 (http://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000276.v2.p1) and described in detail in [21]. The raw data contains 364 590 markers for 5 402 subjects. To obtain roughly independent SNPs this data set was initially screened such that in the final data set the maximal correlation between any pair of SNPs does not exceed $\sqrt{0.1} = 0.316$. The reduced data set contains $p = 26\,233$ SNPs.

The explanatory variables for our genetic model were defined in Table 1, where *a* denotes the less frequent (variant) allele. In case when population frequencies of both alleles are the same, variables \tilde{X} and \tilde{Z} are uncorrelated. In other cases correlations between these variables is different from zero and can be very strong for rare genetic variants. Since each SNP is described by two dummy variables, the full design matrix $[\tilde{X} \ \tilde{Z}]$ contains 52 466 potential regressors. This matrix was then centered and standardized, so the columns of the final design matrix $[X \ Z]$ have zero mean and unit norm.

The trait values are simulated according to two scenarios. In Scenario 1 we simulate from an additive model, where each of the causal SNPs influences the trait only through the additive dummy variable in matrix X ,

$$y = X\beta_X + \epsilon . \tag{2.19}$$

Here $\epsilon \sim \mathcal{N}(0, \mathbf{I})$, the number of ‘causal’ SNPs k varies between 1 and 80 and each causal SNP has an additive effect (non-zero components of β_X) equal to 5 or -5 , with $P(\beta_{X_i} = 5) = P(\beta_{X_i} = -5) = 0.5$. In each of 100 iterations of our experiment causal SNPs were randomly selected from the full set of 26 233 SNPs.

The additive model (2.19) assumes that for each of the SNPs the expected value of the trait for the heterozygote *aA* is the average of expected trait values for both homozygotes *aa* and *AA*. This idealistic assumption is usually not satisfied and many of the SNPs exhibit some dominance effects. To illustrate the performance of gSLOPE in the presence of dominance effects we simulated data according to Scenario 2;

$$y = [X \ Z] \begin{bmatrix} \beta_X \\ \beta_Z \end{bmatrix} + \epsilon \tag{2.20}$$

which differs from Scenario 1 by adding dominance effects (non-zero components of β_Z), which for each of k selected SNPs are sampled from the uniform distribution on $[-5, -3] \cup [3, 5]$. The simulated data sets were analyzed using three different approaches:

- gSLOPE with $p = 26\,233$ groups, where each of the groups contains two explanatory variables, describing the additive and the dominance effect of the same SNP,
- SLOPE $_X$, where the regular SLOPE is used to search through the reduced design matrix X (as in [9] or [11]),
- SLOPE $_{XZ}$, where the regular SLOPE is used to search through the full design matrix $[X \ Z]$.

In all versions of SLOPE we used the iterative procedure for estimation of σ and the sequence λ heuristically adjusted to the case of the Gaussian design matrix, as implemented in the CRAN packages SLOPE and grpsLOPE.

Figure 4 provides the summary of this simulation study. Here FDR and power are calculated at the SNP level. Specifically, in case of SLOPE_{XZ} the SNP is counted as a one discovery if the corresponding additive or the dominance dummy variable is selected.

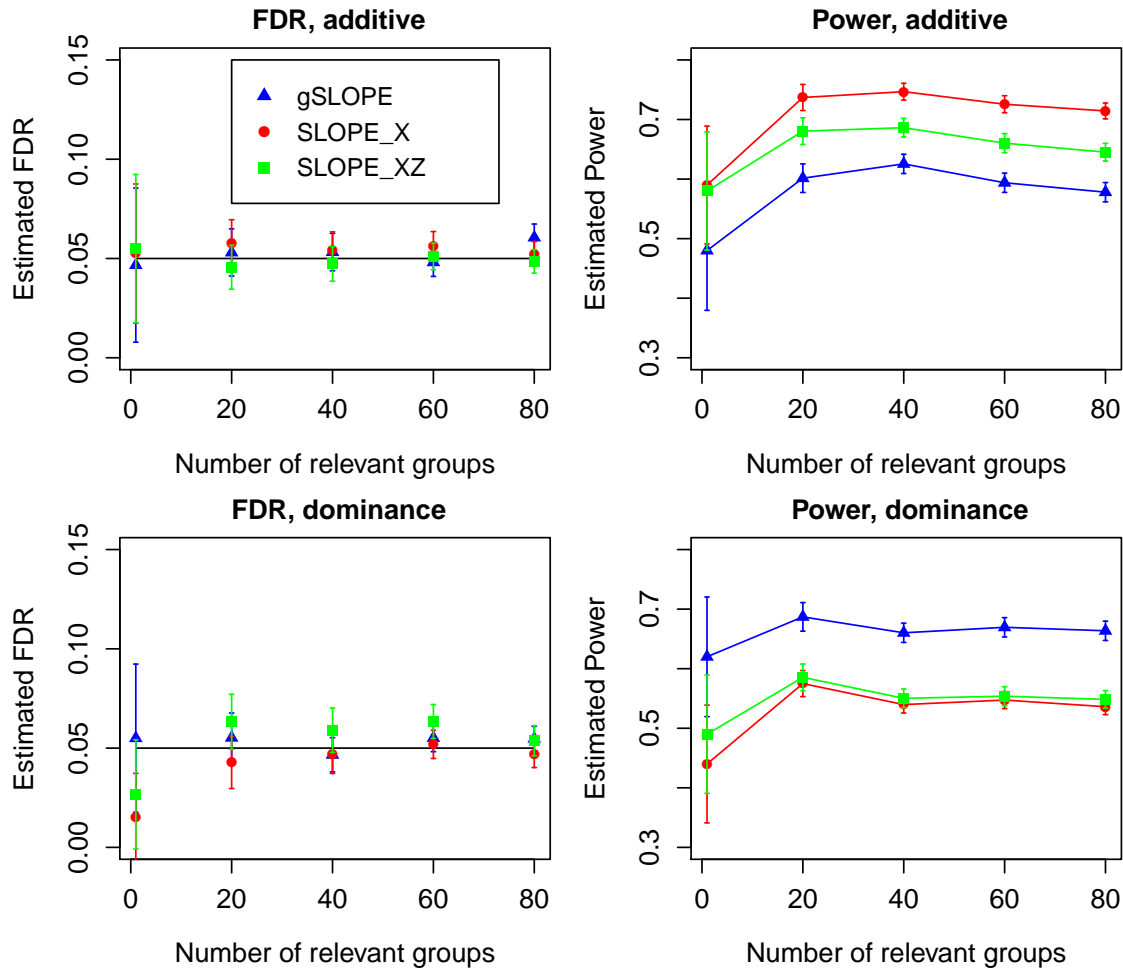


Figure 4: Simulations using real SNP genotypes: $n = 5\,402$, $p = 26\,233$. Power and gFDR are estimated based on 100 iterations of each simulation scenario. Upper panel illustrates the situation where all causal SNPs have only additive effects, while in lower panel each causal SNP has also some dominance effect.

As shown in Figure 4, for both of the simulated scenarios all versions of SLOPE control gFDR for all considered values of k . When the data are simulated according to the additive model the highest power is offered by SLOPE_X , with the power of gSLOPE being smaller by approximately 13% over the whole range of k . However, in the presence of large dominance effects the situation is reversed and gSLOPE offers the highest power, which systematically exceeds the power of SLOPE_X by the symmetric amount of 13%. In our simulations SLOPE_{XZ} has intermediate performance and does not substantially improve the power of SLOPE_X in the presence of dominance effects. Thus our simulations suggest that gSLOPE provides an information complementary to SLOPE_X and might be a useful tool in the context of Genome-Wise Association Studies.

2.9 gSLOPE under GWAS application: real phenotype data

Finally, we have applied group SLOPE to identify SNPs associated with four lipid phenotypes available in NFBC dataset: high-density lipoproteins (HDL), low-density lipoproteins (LDL), triglycerides (TG), and total cholesterol (CHOL). The data set contains genotypes of 364 590 SNPs for 5 402 individuals and was previously analyzed in [11] using regular SLOPE to search for additive SNP effects. Before this analysis the data were reduced by applying the p-value threshold for single marker t-tests and selecting representatives of strongly correlated SNPs, also based on p-values. Since this pre-processing selects most promising SNPs by performing multiple testing on the full set of $p = 364590$ SNPs, the sequence of the tuning parameters for SLOPE needs to be adjusted to this value of p rather than to the number of selected representatives. The algorithm for this analysis is implemented in R package `geneSLOPE` and its details are explained in [11]. According to extensive simulation study and real data analysis reported in [11], `geneSLOPE` allows to control FDR for the analysis with full size GWAS data.

In our data analysis we used three versions of SLOPE: `geneSLOPE` for additive effects (as in [11]), `geneSLOPEXZ`, with the design matrix extended by inclusion of dominance dummy variables, and `gene group SLOPE` (`geneGSLOPE`). In `geneSLOPEXZ` and `geneGSLOPE` representative SNPs were selected based on the one way ANOVA tests. For all these procedures the pre-processing was based on p-value threshold $p < 0.05$ and the correlation cutoff $\rho < 0.3$, which allowed to reduce the data set to roughly 8500 of interesting representative SNPs. For the convenience of the reader, the Procedure 3 for the full `geneGSLOPE` analysis is provided below.

Procedure 3 `geneGSLOPE` procedure

Input: $r \in (0, 1)$, $\pi \in (0, 1]$

Screen SNPs:

- (1) For each SNP calculate independently the p -value for the ANOVA test with the null hypothesis, $H_0 : \mu_{aa} = \mu_{aA} = \mu_{AA}$.
- (2) Define the set \mathcal{B} of indices corresponding to SNPs whose p -values are smaller than π .

Cluster SNPs:

- (3) Select the SNP j in \mathcal{B} with the smallest p -value and find all SNPs whose Pearson correlation with this selected SNP is larger than or equal to r .
- (4) Define this group as a cluster and SNP j as the representative of the cluster. Include SNP j in \mathcal{S} , and remove the entire cluster from \mathcal{B} .
- (5) Repeat steps (3)-(4) until \mathcal{B} is empty. Denote by m number of all clumps (this is also the number of elements in \mathcal{S}).

Selection:

- (6) Apply the iterative gSLOPE method (i.e. gSLOPE with σ estimation and correction for independent regressors) on $X_{\mathcal{S}}$, being matrix X restricted to columns corresponding to the set \mathcal{S} of selected SNPs. Here, the tuning parameters, vector λ , is defined as in Procedure 1, with p being the number of all initial SNPs, and then this vector is restricted only to first m coefficients.
 - (7) Representatives which were selected indicate the selection of entire clumps.
-

Results in the context of number of discoveries given by `geneSLOPE`, `geneSLOPEXZ` and `geneGSLOPE` are summarized in Table 2, where we can observe that both `geneSLOPE` and `geneSLOPEXZ`, gave identical results for LDL, CHOL and TG. Compared to these methods `geneGSLOPE` did not reveal any new response-related SNPs for LDL and CHOL. Actually, for these two traits `geneGSLOPE` missed some SNPs detected by the other two methods.

A different situation takes place for TG, where `geneGSLOPE` identifies 6 additional SNPs as compared to the other two methods. All these detections have a similar structure, showing a significant recessive effect of the minor allele. In all these cases the minor allele frequency was smaller than 0.1. The detection of such "rare" recessive effects by the simple linear regression

	HDL	LDL	TG	CHOL
geneSLOPE	7	6	2	5
geneSLOPE _{XZ}	8	6	2	5
geneGSLOPE	8	4	8	4
New discoveries: geneSLOPE _{XZ}	1	0	0	0
New discoveries: geneGSLOPE	2	0	6	0

Table 2: Number of discoveries in real data analysis

model is rather difficult, since the regression line adjusts mainly to the two prevalent genotype groups and is almost flat [19].

In case of HDL all three versions of SLOPE gave different results. geneSLOPE_{XZ} identifies one new SNP as compared to geneSLOPE, while geneGSLOPE identifies one more SNP and misses one of the discoveries obtained by other two methods. In Figure 5 we compare two exemplary discoveries: one detected at the same time by geneSLOPE and geneGSLOPE (known discovery) and one detected only by geneGSLOPE (new discovery). This example clearly shows the additive effect of the previously detected SNP and the recessive character of the second SNP. In case of new discovery there are only 5 individuals in the last genotype group, which makes the change in the mean not detectable by simple linear regression.

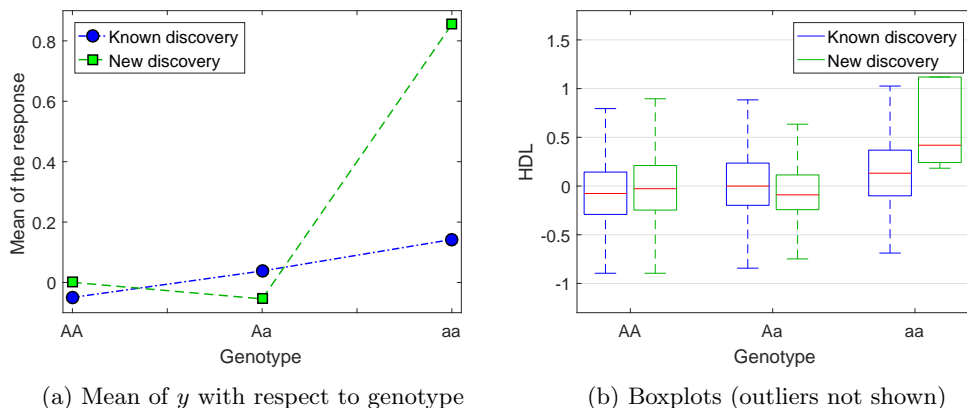


Figure 5: Comparison of discovery detected by both geneSLOPE and geneGSLOPE (known discovery), and discovery detected only by geneGSLOPE (new discovery).

The results of real data analysis agree with results of simulations. They show that geneGSLOPE has a lower power than geneSLOPE for detection of additive effects but can be very helpful in detecting rare recessive variants. Thus these two methods are complementary to each other and can be used together to enhance the power of detection of influential genes.

3 Discussion

Group SLOPE is a new convex optimization procedure for selection of important groups of explanatory variables, which can be considered as a generalization of group LASSO and of SLOPE. In this article we provide an algorithm for solving group SLOPE and discuss the choice of the sequence of regularizing parameters. Our major focus is the control of group FDR, which can be obtained when variables in different groups are orthogonal to each other or they are stochastically inde-

pendent and the signal is sufficiently sparse. After some preprocessing of the data such situations occur frequently in the context of genetic studies, which in this paper serve as a major example of applications. While we concentrated mainly on using gSLOPE to group dummy variables corresponding to different effects of the same SNP, gSLOPE can be used to group SNPs based on biological function, physical location etc. We also expect this method to be advantageous in the context of identification of groups of rare genetic variants, where considering their joint effect on phenotype should substantially increase the power of detection.

The major purpose of controlling FDR rather than absolutely eliminating false discoveries is the wish to increase the power of detection of signals which are comparable to the noise level. As shown by a variety of theoretical and empirical results, this allows SLOPE to obtain an optimal balance between the number of false and true discoveries and leads to very good estimation and predictive properties (see e.g. [10], [9] or [26]). Our Theorem 2.6 illustrates that these good estimation properties are inherited by group SLOPE.

We provide the regularizing sequence λ^{max} , which provably controls gFDR in case when variables in different groups are orthogonal. Additionally, we propose its relaxation λ^{mean} , which according to our extensive simulations controls "average" gFDR, where the average is with respect to all possible signal placements. This sequence can be easily modified taking into account the prior distribution on the signal placement. Such "Bayesian" version of gSLOPE and the proof of control of the respective average gFDR remains an interesting topic for a further research.

Another important topic for a further research is the formal proof of gFDR control when variables in different groups are independent and setting precise limits on the sparsity levels under which it can be done. Asymptotic formulas, which allow for very accurate prediction of FDR for LASSO under Gaussian design are provided in [25]. We expect that similar results can be obtained for SLOPE and gSLOPE and generalized to the case of random matrices, where variables are independent and come from sub-Gaussian distributions. However, the technical complexity of results reported in [25] illustrates that this task is rather challenging. An alternative approach for the perfect gFDR control under random designs is to couple gSLOPE with the new knock-off procedure proposed in [12]. We expect that such a combination should allow to increase the power of detection of relevant features, as compared to other methods currently used with knock-offs.

While we concentrated on control of FDR in case when groups of variables are roughly orthogonal to each other, it is worth mentioning that original SLOPE has very interesting properties also in case when regressors are strongly correlated. As shown e.g. in [14], the Sorted L-One norm has a tendency to average estimated regression coefficients over groups of strongly correlated predictors, which enhances the predictive properties. This also allows not to lose important predictors due to their correlation with other features. We expect similar properties to hold for gSLOPE but the investigation of the properties of gSLOPE when variables in different groups are strongly correlated remains an interesting topic for a further research.

Acknowledgement

We would like to thank Emmanuel J. Candès and Jan Mielniczuk for helpful remarks and suggestions and Christine Peterson for screening the North Finland Birth Cohort (NFBC) dataset. D. B. would like to thank Professor Jerzy Ombach for significant help with the process of obtaining access to the data. D. B. and M. B. are supported by European Union's 7th Framework Programme for research, technological development and demonstration under Grant Agreement no 602552 and by the Polish Ministry of Science and Higher Education according to agreement 2932/7.PR/2013/2. Additionally D.B. acknowledges the support from NIMH grant R01MH108467, W. S. was partially supported by a General Wang Yaowu Stanford Graduate Fellowship.

References

- [1] Abramovich F. and Benjamini Y. and Donoho D. L. and Johnstone I. M. Adapting to unknown sparsity by controlling the false discovery rate. *Ann. Statist.*, 34(2):584–653, 2006.
- [2] Abramowitz M. and Stegun I. A. *Handbook of mathematical functions: with formulas, graphs, and mathematical tables*. Number 55. Courier Corporation, 1964.
- [3] Akaike H. A New Look at the Statistical Model Identification. *IEEE Transactions on Automatic Control*, 19(6):716–723, 1974.
- [4] Bakin S. Adaptive regression and model selection in data mining problems. 1999.
- [5] Barber R. F. and Candès E. J. Controlling the False Discovery Rate via Knockoffs. *The Annals of Statistics*, 43(5):2055–2085, 2015.
- [6] Beck A. and Teboulle M. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 1:183–202, 2009.
- [7] Benjamini Y. and Hochberg Y. Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society, Series B*, 57(1):289–300, 1995.
- [8] Bogdan M. and Chakrabarti A. and Frommlet F. and Ghosh J. K. Asymptotic Bayes Optimality under sparsity of some multiple testing procedures. *Annals of Statistics*, 39:1551–1579, 2011.
- [9] Bogdan M. and van den Berg E. and Sabatti C. and Su W. and Candès E. J. SLOPE – adaptive variable selection via convex optimization. *Annals of Applied Statistics*, 9(3):1103–1140, 2015.
- [10] Bogdan M. and van den Berg E. and Su W. and Candès E. J. Statistical Estimation and Testing via the Ordered ℓ_1 Norm. *arXiv:1310.1969*, 2013.
- [11] Brzyski, D. and Peterson, C.B. and Sobczyk, P. and Candès, E.J. and Bogdan, M. and Sabatti, C. Controlling the rate of GWAS false discoveries. *bioRxiv 058230*, to appear in *Genetics*, 2016.
- [12] Candès, E.J and Fan, Y. and Janson L. and Lv J. Panning for gold: model-free knockoffs for high-dimensional controlled variable selection. *arXiv:1610.02351*, 2016.
- [13] Donoho D. L. and Johnstone I. M. Minimax risk over ℓ_p -balls for ℓ_q -error. *Probability Theory and Related Fields*, 99(2):277–303, 1994.
- [14] Figueiredo M. A. T. and Nowak R. D. Ordered Weighted l_1 Regularized Regression with Strongly Correlated Covariates: Theoretical Aspects. *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics, JMLR:W&CP*, 51:930–938, 2016.
- [15] Frommlet F. and Bogdan M. Some optimality properties of FDR controlling rules under sparsity. *Electronic Journal of Statistics*, 7:1328–1368, 2013.
- [16] Gossmann A. and Cao S. and Wang Y.-P. Identification of significant genetic variants via SLOPE”, and its extension to Group SLOPE. In *Proceedings of the International Conference on Bioinformatics, Computational Biology and Biomedical Informatics*, 2015.
- [17] Hardy G.H. and Littlewood, J.E. and Pólya, G. *Inequalities*. 1952.
- [18] Inglot T. Inequalities for quantiles of the chi-square distribution. *Probability and Mathematical Statistics*, 30(2):339–351, 2010.

- [19] Lettre G. and Lange C. and Hirschhorn J. N. Genetic model testing and statistical power in population-based association studies of quantitative traits. *Genetic Epidemiology*, 31(4): 358–362, 2007.
- [20] Neuvial P. and Roquain E. On false discovery rate thresholding for classification under sparsity. *Annals of Statistics*, 40:2572–2600, 2012.
- [21] Sabatti C. and Service S. K. and Hartikainen A. and Pouta A. and Ripatti S. and Brodsky J. and Jones C. G. and Zaitlen N. A. and Varilo T. and Kaakinen M. and Sovio U. and Ruokonen A. and Laitinen J. and Jakkula E. and Coin L. and Hoggart C. and Collins A. and Turunen H. and Gabriel S. and Elliot P. and McCarthy M. I. and Daly M. J. and Jvelin M. and Freimer N. B. and Peltonen L. Genome-wide association analysis of metabolic traits in a birth cohort from a founder population. *Nature Genetics*, 41(1):35–46, 2009.
- [22] Schwarz G. Estimating the Dimension of a Model. *The Annals of Statistics*, 6(2):461–464, 1978.
- [23] Simon N. and Friedman J. and Hastie T. and Tibshirani R. A sparse-group lasso. *Journal of Computational and Graphical Statistics*, 2013.
- [24] Simon N. and Tibshirani R. Standardization and the group lasso penalty. Technical report, 2011.
- [25] Su W. and Bogdan M. and Candès E.J. False discoveries occur early on the lasso path. *arXiv:1511.01957*, to appear in *Ann. Statist.*, 2015.
- [26] Su W. and Candès E. SLOPE is adaptive to unknown sparsity and asymptotically minimax. *Annals of Statistics*, 40:1038–1068, 2016.
- [27] Tibshirani R. Regression Shrinkage and Selection Via the Lasso. *Journal of the Royal Statistical Society, Series B*, 58:267–288, 1996.
- [28] Tseng P. On accelerated proximal gradient methods for convex-concave optimization. 2008.
- [29] Yuan M. and Lin Y. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society, Series B*, 68(1):49–67, 2006.

A J_λ norm properties

For nonnegative, nonincreasing sequence $\lambda_1 \geq \dots \geq \lambda_p \geq 0$ consider function $\mathbb{R}^p \ni b \mapsto J_\lambda(b) \in \mathbb{R}$ given by $J_\lambda(b) = \sum_{i=1}^p \lambda_i \cdot |b|_{(i)}$, where $|b|_{(1)} \geq \dots \geq |b|_{(p)}$ is the vector of sorted absolute values.

Proposition A.1. *If $a, b \in \mathbb{R}^p$ are such that $|a| \preceq |b|$, then $|a|_{(\cdot)} \preceq |b|_{(\cdot)}$.*

Proof. Without loss of generality we can assume that a and b are nonnegative and that it occurs $a_1 \geq \dots \geq a_p$. We will show that $a_k \leq b_{(k)}$ for $k \in \{1, \dots, p\}$. Fix such k and consider the set $S_k := \{b_i : b_i \geq a_k\}$. It is enough to show that $|S_k| \geq k$. For each $j \in \{1, \dots, k\}$ we have

$$b_j \geq a_j \geq a_k \implies b_j \in S_k,$$

what proves the last statement. ■

Corollary A.2. *Let $a \in \mathbb{R}^p, b \in \mathbb{R}^p$ and $|a| \preceq |b|$ then Proposition (A.1) instantly gives that $J_\lambda(a) \leq J_\lambda(b)$, since $J_\lambda(a) = \lambda^\top |a|_{(\cdot)} \leq \lambda^\top |b|_{(\cdot)} = J_\lambda(b)$.*

Proposition A.3. *For fixed sequence $\lambda_1 \geq \dots \geq \lambda_p \geq 0$, let $b \in \mathbb{R}^p$ be such that $b \succeq 0$ and $b_j > b_l$ for some $j, l \in \{1, \dots, p\}$. For $0 < \varepsilon \leq (b_j - b_l)/2$, define $b_\varepsilon \in \mathbb{R}^p$ by conditions $(b_\varepsilon)_l := b_l + \varepsilon$, $(b_\varepsilon)_j := b_j - \varepsilon$ and $(b_\varepsilon)_i := b_i$ for $i \notin \{j, l\}$. Then $J_\lambda(b_\varepsilon) \leq J_\lambda(b)$.*

Proof. Let $\pi : \{1, \dots, p\} \rightarrow \{1, \dots, p\}$ be permutation such as $\sum_{i=1}^p \lambda_i (b_\varepsilon)_{(\pi(i))} = \sum_{i=1}^p \lambda_{\pi(i)} (b_\varepsilon)_i$ for each i in $\{1, \dots, p\}$ and $\lambda_{\pi(j)} \geq \lambda_{\pi(l)}$. From the rearrangement inequality (Theorem 368 in [17]),

$$\begin{aligned} J_\lambda(b) - J_\lambda(b_\varepsilon) &= \sum_{i=1}^p \lambda_i b_{(i)} - \sum_{i=1}^p \lambda_i (b_\varepsilon)_{(i)} = \sum_{i=1}^p \lambda_i b_{(i)} - \sum_{i=1}^p \lambda_{\pi(i)} (b_\varepsilon)_i \\ &\geq \sum_{i=1}^p \lambda_{\pi(i)} b_i - \sum_{i=1}^p \lambda_{\pi(i)} (b_\varepsilon)_i = \varepsilon (\lambda_{\pi(j)} - \lambda_{\pi(l)}) \geq 0. \end{aligned} \tag{A.1}$$

■

B Numerical algorithm

In this section we will discuss the convexity of the objective function and the algorithm for computing the solution to gSLOPE problem (2.2). Our optimization method is based on the fast algorithm for evaluation of the proximity operator (prox) for sorted ℓ_1 norm, which was derived in [10].

B.1 Convexity of the objective function

To show that the objectives in problems (2.2) and (2.5) are convex functions, we will prove the following propositions

Proposition B.1. *Function $J_{\lambda, W, \mathbb{I}}(b) := J_\lambda(W \llbracket b \rrbracket_{\mathbb{I}})$ is a norm for any nonnegative, nonincreasing sequence $\{\lambda_i\}_{i=1}^m$ containing at least one nonzero element, partition \mathbb{I} of the set $\{1, \dots, \tilde{p}\}$ and diagonal matrix W with positive elements on diagonal.*

Proof. It is easy to see that $J_{\lambda, W, \mathbb{I}}(c) = 0$ if and only if $c = 0$ and that for any scalar $\alpha \in \mathbb{R}$ it occurs $J_{\lambda, W, \mathbb{I}}(\alpha c) = |\alpha| J_{\lambda, W, \mathbb{I}}(c)$. We will show that $J_{\lambda, W, \mathbb{I}}$ satisfies the triangle inequality. Let b, c be any vectors from $\mathbb{R}^{\tilde{p}}$. From the positivity of w_i 's we have $W \llbracket a + b \rrbracket_{\mathbb{I}} \preceq W \llbracket a \rrbracket_{\mathbb{I}} + W \llbracket b \rrbracket_{\mathbb{I}}$. Therefore, Corollary A.2 yields

$$\begin{aligned} J_{\lambda, W, \mathbb{I}}(a + b) &= J_\lambda(W \llbracket a + b \rrbracket_{\mathbb{I}}) \leq J_\lambda(W \llbracket a \rrbracket_{\mathbb{I}} + W \llbracket b \rrbracket_{\mathbb{I}}) \\ &\leq J_\lambda(W \llbracket a \rrbracket_{\mathbb{I}}) + J_\lambda(W \llbracket b \rrbracket_{\mathbb{I}}) = J_{\lambda, W, \mathbb{I}}(a) + J_{\lambda, W, \mathbb{I}}(b), \end{aligned} \tag{B.1}$$

since J_λ is a norm. ■

Proposition B.2. *Function $J_\lambda(W\llbracket b\rrbracket_{X,I})$ is a seminorm for any nonnegative, nonincreasing sequence $\{\lambda_i\}_{i=1}^m$, partition I of the set $\{1, \dots, p\}$, design matrix $X \in M(n, p)$ and diagonal matrix W with positive elements on diagonal.*

Proof. Clearly, $J_\lambda(W\llbracket \alpha b\rrbracket_{X,I}) = |\alpha|J_\lambda(W\llbracket b\rrbracket_{X,I})$, for any scalar $\alpha \in \mathbb{R}$. Moreover, for any $a, b \in \mathbb{R}^p$, it holds $W\llbracket a + b\rrbracket_{X,I} \preceq W\llbracket a\rrbracket_{X,I} + W\llbracket b\rrbracket_{X,I}$, and the triangle inequality could be proved similarly as in the previous proposition. \blacksquare

B.2 Proximal gradient method

Consider unconstrained optimization problem of form

$$\underset{b}{\text{minimize}} \quad f(b) = g(b) + h(b), \quad (\text{B.2})$$

where g and h are convex functions and g is differentiable (for example LASSO and SLOPE are of such form). There exist efficient methods, namely *proximal gradient algorithms*, which could be applied to find numerical solution for such objective functions. To design efficient algorithms, however, h must be prox-capable, meaning that there is known fast algorithm for computing the proximal operator for h ,

$$\text{prox}_{th}(y) := \arg \min_b \left\{ \frac{1}{2t} \|y - b\|_2^2 + h(b) \right\}, \quad (\text{B.3})$$

for each $y \in \mathbb{R}^p$ and $t > 0$. The iterative algorithm works as follows. Suppose that in k step $b^{(k)}$ is the current guess. Then, guess $b^{(k+1)}$ is given by

$$b^{(k+1)} := \arg \min_b \left\{ g(b^{(k)}) + \langle \nabla g(b^{(k)}), b - b^{(k)} \rangle + \frac{1}{2t} \|b - b^{(k)}\|_2^2 + h(b) \right\}. \quad (\text{B.4})$$

The two first terms in objective function in (B.4) are Taylor approximation of g , third addend is a proximity term which is responsible for searching an update reasonably close and t can be treated as a step size.

Problem (B.4) could be reformulated to

$$b^{(k+1)} := \arg \min_b \left\{ \frac{1}{2} \|b^{(k)} - t\nabla g(b^{(k)}) - b\|_2^2 + th(b) \right\}, \quad (\text{B.5})$$

hence $b^{(k+1)} = \text{prox}_{th}(b^{(k)} - t\nabla g(b^{(k)}))$, which justifies the need for existence of a fast algorithm computing values of the proximal operator. In each step the value of t could be changed raising the sequence $\{t_i\}_{i=1}^\infty$. In situation when $g(b) = \frac{1}{2} \|y - Xb\|_2^2$, we get the following algorithm.

Procedure 4 Proximal gradient algorithm

input: $b^{[0]} \in \mathbb{R}^p$, $k=0$
while (Stopping criteria are not satisfied) **do**
 1. $b^{[k+1]} = \text{prox}_{t_k h_\lambda}(b^{[k]} - t_k X^\top (Xb^{[k]} - y))$;
 2. $k \leftarrow k + 1$.
end while

It is known that t_i 's could be selected in different ways to ensure that $f(b^{(k)})$ converges to the optimal value [6], [28].

B.3 Proximal operator for gSLOPE

Let $I = \{I_1, \dots, I_m\}$, l_i be rank of submatrix X_{I_i} for $i = 1, \dots, m$ and $\lambda = (\lambda_1, \dots, \lambda_m)^\top$ be a vector satisfying $\lambda_1 \geq \dots \geq \lambda_m \geq 0$. We will now employ the proximal gradient method to find the numerical solution to (2.2). As stated in subsection 2.1, we can focus on the equivalent optimization problem (2.5), namely we aim to solve problem

$$b^* := \arg \min_b \left\{ \frac{1}{2} \|y - \tilde{X}b\|_2^2 + \sigma J_\lambda(W \llbracket b \rrbracket_{\mathbb{I}}) \right\}, \quad (\text{B.6})$$

with $\mathbb{I} = \{\mathbb{I}_1, \dots, \mathbb{I}_m\}$ being a partition of the set $\{1, \dots, \tilde{p}\}$ where $\tilde{p} = l_1 + \dots + l_m$.

Without loss of generality we assume that $\sigma = 1$. Since considered objective is of form (B.2), we can apply proximal gradient algorithm, provided that norm $J_{\lambda, \mathbb{I}, W}$ is prox-capable. To compute the proximal operator for $J_{\lambda, \mathbb{I}, W}$ we must be able to minimize $\frac{1}{2t} \|y - b\|_2^2 + J_{\lambda, \mathbb{I}, W}(b)$, for any $y \in \mathbb{R}^{\tilde{p}}$ and $t > 0$. Multiplying objective by positive number, t , does not change the solution. Such operation leads to new objective function, $\frac{1}{2} \|y - b\|_2^2 + J_{t\lambda, \mathbb{I}, W}(b)$. This shows that it is enough to derive a fast algorithm for finding the numerical solution to the problem

$$\text{prox}_J(y) := \arg \min_b \left\{ \frac{1}{2} \|y - b\|_2^2 + J_{\lambda, \mathbb{I}, W}(b) \right\}, \quad (\text{B.7})$$

which could be applicable to arbitrary sequence $\lambda_1 \geq \dots \geq \lambda_m \geq 0$.

We will start with situation when W is identity matrix. Simply, then $\text{prox}_J(y)$ is proximal operator for function $J_{\lambda, \mathbb{I}}(b) := J_\lambda(\llbracket b \rrbracket_J)$. In such a case computing (B.7) could be immediately reduced to finding prox for J_λ norm, since thanks to (2.8) we have

$$\begin{cases} c^* = \arg \min_c \left\{ \frac{1}{2} \|\llbracket y \rrbracket_{\mathbb{I}} - c\|_2^2 + J_\lambda(c) \right\} \\ (\text{prox}_J(y))_{\mathbb{I}_i} = c_i^* (\|y_{\mathbb{I}_i}\|_2)^{-1} y_{\mathbb{I}_i}, \quad i = 1, \dots, m \end{cases}. \quad (\text{B.8})$$

Consequently, $\text{prox}_J(y)$ could be obtained by applying two steps procedure: find c^* by using fast prox algorithm for J_λ for vector $\llbracket y \rrbracket_{\mathbb{I}}$, and compute $\text{prox}_J(y)$ by applying simple calculus to c^* .

Consider now general situation with fixed positive numbers w_1, \dots, w_m and define diagonal matrix M by conditions $M_{\mathbb{I}_i, \mathbb{I}_i} := w_i^{-1} \mathbf{I}_{l_i}$, for $i = 1, \dots, m$. Then

$$J_{\lambda, \mathbb{I}, W}(b) = J_\lambda(W \llbracket b \rrbracket_{\mathbb{I}}) = J_\lambda(\llbracket M^{-1}b \rrbracket_{\mathbb{I}}) = J_{\lambda, \mathbb{I}}(M^{-1}b). \quad (\text{B.9})$$

Since M is nonsingular, we can substitute $\eta := M^{-1}b$ and consider equivalent formulation of (B.6)

$$\begin{cases} \eta^* := \arg \min_{\eta} \left\{ \frac{1}{2} \|y - \tilde{X}M\eta\|_2^2 + J_{\sigma\lambda, \mathbb{I}}(\eta) \right\}, \\ b^* = M\eta^* \end{cases}. \quad (\text{B.10})$$

Therefore, after modifying the design matrix, gSLOPE can be always recast as problem with unit weights. Since $J_{\lambda, \mathbb{I}}$ is prox-capable, applying proximal gradient method to (B.10) is straightforward. To implement the method introduced in this article, we have used a modified version of Procedure 4, the accelerated proximal gradient method known as FISTA [6]. In particular FISTA gives a precise procedure for choosing steps sizes, to achieve a fast convergence rate. To derive proper stopping criteria, we have considered dual problem to gSLOPE, described in the following section, and employed the strong duality property.

B.4 Dual norm and conjugate of grouped sorted ℓ_1 norm

Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be a norm. We will use notation f^D to refer to the dual norm to f , i.e function defined as $f^D(x) := \max_b \{x^\top b : f(b) \leq 1\}$. It could be shown (see [9]), that the set C_λ , defined

as $C_\lambda := \left\{ x \in \mathbb{R}^p : \sum_{i=1}^k |x|_{(i)} \leq \sum_{i=1}^k \lambda_i, k = 1, \dots, p \right\}$, is unit ball of the dual norm to J_λ for any nonnegative, nonincreasing sequence $\{\lambda_i\}_{i=1}^p$ with at least one nonzero element. We will now consider the dual norm to $J_{\lambda,I,W}(b) = J_\lambda(W\llbracket b \rrbracket_I)$. It holds

$$\begin{aligned} J_{\lambda,I,W}^D(x) &= \max_b \{x^\top b : J_{\lambda,I,W}(b) \leq 1\} = \max_b \{x^\top b : J_\lambda(W\llbracket b \rrbracket_I) \leq 1\} = \\ &= \max_{b,c} \{x^\top b : J_\lambda(c) \leq 1, c = W\llbracket b \rrbracket_I\} = \max_c \{x^\top b^c : J_\lambda(c) \leq 1, c \succeq 0\}, \end{aligned} \quad (\text{B.11})$$

where b^c is defined as $b^c := \arg \max_b \{x^\top b : c = W\llbracket b \rrbracket_I\}$. This problem is separable and for each i we have $b_{I_i}^c = \arg \max \{x_{I_i}^\top b_{I_i} : c_i^2 = w_i^2 \|b_{I_i}\|_2^2\}$. Applying the Lagrange multiplier method quickly yields $x_{I_i}^\top b_{I_i}^c = c_i w_i^{-1} \|x_{I_i}\|_2$. Consequently,

$$\begin{aligned} J_{\lambda,I,W}^D(x) &= \max_c \{(W^{-1}\llbracket x \rrbracket_I)^\top c : J_\lambda(c) \leq 1, c \succeq 0\} = \\ &= \max_c \{(W^{-1}\llbracket x \rrbracket_I)^\top c : J_\lambda(c) \leq 1\} = J_\lambda^D(W^{-1}\llbracket x \rrbracket_I). \end{aligned} \quad (\text{B.12})$$

Therefore, $\{x : J_{\lambda,I,W}^D(x) \leq 1\} = \{x : J_\lambda^D(W^{-1}\llbracket x \rrbracket_I) \leq 1\} = \{x : W^{-1}\llbracket x \rrbracket_I \in C_\lambda\}$. Since the conjugate of norm is equal to zero for arguments from unit ball of dual norm, and equal to infinity otherwise, we immediately get

Corollary B.3. *The conjugate function for $J_{\lambda,I,W}$ is the function $J_{\lambda,I,W}^*$ defined as*

$$J_{\lambda,I,W}^*(x) = \begin{cases} 0, & W^{-1}\llbracket x \rrbracket_I \in C_\lambda \\ \infty, & \text{otherwise} \end{cases}. \quad (\text{B.13})$$

B.5 Stopping criteria for numerical algorithm

Without loss of generality assume that $\sigma = 1$. We will start with optimization problem in (B.10), namely

$$\underset{\eta}{\text{minimize}} \quad f(\eta) = \frac{1}{2} \|y - \tilde{X}M\eta\|_2^2 + J_\lambda(\llbracket \eta \rrbracket_{\mathbb{I}}) \quad (\text{B.14})$$

for $\llbracket \eta \rrbracket_{\mathbb{I}} = (\|\eta_{\mathbb{I}_1}\|_2, \dots, \|\eta_{\mathbb{I}_m}\|_2)^\top$ and $M_{\mathbb{I}_i, \mathbb{I}_i} = \frac{1}{w_i} \mathbf{I}_i$, $i = 1, \dots, m$. This problem could be written in equivalent form

$$\begin{aligned} \underset{\eta, r, c}{\text{minimize}} \quad & \frac{1}{2} \|r\|_2^2 + c \\ \text{s.t.} \quad & \begin{cases} J_{\lambda, \mathbb{I}}(\eta) - c \leq 0 \\ y - r - \tilde{X}M\eta = 0 \end{cases} \end{aligned} \quad (\text{B.15})$$

(notice that for (η^*, r^*, c^*) being solution, it must occurs $c^* = J_{\lambda, \mathbb{I}}(\eta^*)$). Since (B.15) is convex and (η_0, r_0, c_0) , for $\eta_0 = 0$, $r_0 = y$ and $c_0 = 1$, is strictly feasible, the strong duality holds. Lagrange dual function for this problem is given by

$$\begin{aligned} g(\mu, \nu) &= \inf_{\eta, r, c} \left\{ \frac{1}{2} \|r\|_2^2 + c + \mu^\top (y - r - \tilde{X}M\eta) + \nu (J_{\lambda, \mathbb{I}}(\eta) - c) \right\} = \\ &= \mu^\top y + \inf_r \left\{ \frac{1}{2} \|r\|_2^2 - \mu^\top r \right\} + \inf_c \{c - \nu c\} + \inf_\eta \{ -\mu^\top \tilde{X}M\eta + \nu J_{\lambda, \mathbb{I}}(\eta) \}. \end{aligned} \quad (\text{B.16})$$

Now, since the minimum of $\frac{1}{2} \|r\|_2^2 - \mu^\top r$ is taken for $r = \mu$, we have

$$g(\mu, \nu) = \mu^\top y - \frac{1}{2} \|\mu\|_2^2 + \inf_c \{c - \nu c\} - J_{\nu\lambda, \mathbb{I}}^*((\tilde{X}M)^\top \mu). \quad (\text{B.17})$$

Then $\nu^* = 1$ and from Corollary B.3, the dual problem to (B.15) is equivalent to

$$\begin{aligned} & \underset{\mu}{\text{maximize}} && \mu^\top y - \frac{1}{2} \|\mu\|_2^2 \\ & \text{s.t.} && \llbracket M\tilde{X}^\top \mu \rrbracket_{\mathbb{I}} \in C_\lambda \end{aligned} \quad (\text{B.18})$$

Let (η^*, r^*, c^*) be primal and $(\mu^*, 1)$ be dual solution to (B.15). Obviously, $\mu^* = r^* = y - \tilde{X}M\eta^*$ and $c^* = J_{\lambda, \mathbb{I}}(\eta^*)$. Furthermore, from strong duality we have

$$\frac{1}{2} \|y - \tilde{X}M\eta^*\|_2^2 + J_{\lambda, \mathbb{I}}(\eta^*) = (y - \tilde{X}M\eta^*)^\top y - \frac{1}{2} \|y - \tilde{X}M\eta^*\|_2^2, \quad (\text{B.19})$$

which gives $(\tilde{X}M\eta^*)^\top (y - \tilde{X}M\eta^*) = J_{\lambda, \mathbb{I}}(\eta^*)$. Now, for current approximate $\eta^{[k]}$ of solution to (B.14), achieved after applying proximal gradient method, we define the current duality gap for k step as

$$\rho(\eta^{[k]}) = (\tilde{X}M\eta^{[k]})^\top (y - \tilde{X}M\eta^{[k]}) - J_{\lambda, \mathbb{I}}(\eta^{[k]}) \quad (\text{B.20})$$

and we will determine the infeasibility of $\mu^{[k]} := y - \tilde{X}M\eta^{[k]}$ by using the measure

$$\text{infeas}(\mu^{[k]}) := \max \left\{ J_{\lambda, \mathbb{I}}^D(M\tilde{X}^\top \mu^{[k]}) - 1, 0 \right\} \quad (\text{B.21})$$

To define the stopping criteria we have applied the widely used procedure: treat $\rho(\eta^{[k]})$ as indicator telling how far $\eta^{[k]}$ is from true solution and terminate the algorithm when this difference and infeasibility measure are sufficiently small. Summarizing, we have derived algorithm according to scheme

Procedure 5 group SLOPE

input: infeas.tol: *positive number determining the tolerance for infeasibility*;
dual.tol: *positive number determining the tolerance for duality gap*;
 $k := 0$, $\eta^{[0]}$, $\mu^{[0]} := \mu(\eta^{[0]})$, $\text{infeas}^{[0]} := \text{infeas}(\mu^{[0]})$, $\rho^{[0]} := \rho(\eta^{[0]})$;
while ($\text{infeas}^{[k]} > \text{infeas.tol}$ or $\rho^{[k]} > \text{dual.tol}$) **do**
1. $k \leftarrow k + 1$;
2. get $\eta^{[k]}$ from Procedure 4;
3. $\mu^{[k]} := \mu(\eta^{[k]})$;
4. $\text{infeas}^{[k]} := \text{infeas}(\mu^{[k]})$, $\rho^{[k]} := \rho(\eta^{[k]})$;
end while
 $\beta_{gS} := M\eta^{[k]}$.

C Alternative representation in the orthogonal case

Suppose that the experiment matrix is orthogonal at group level, i.e. it holds $X_{I_i}^\top X_{I_j} = \mathbf{0}$, for every $i, j \in \{1, \dots, m\}$, $i \neq j$. In such a case, \tilde{X} in problem (2.5) is orthogonal matrix, i.e. $\tilde{X}^\top \tilde{X} = \mathbf{I}_{\tilde{p}}$. If $n = \tilde{p}$, i.e. \tilde{X} is a square and orthogonal matrix, we also have $\tilde{X}\tilde{X}^\top = \mathbf{I}_{\tilde{p}}$ and it obeys $\|\tilde{X}^\top b\|_2^2 = b^\top \tilde{X}\tilde{X}^\top b = \|b\|_2^2$ for $b \in \mathbb{R}^{\tilde{p}}$. For the general case with $n \geq \tilde{p}$, we can extend \tilde{X} to a square matrix by adding new orthonormal columns and defining $\tilde{X}_C := [\tilde{X} \ C]$, where C is composed of vectors (columns) being some complement to orthogonal basis of $\mathbb{R}^{\tilde{p}}$. For $y \in \mathbb{R}^n$ and $b \in \mathbb{R}^{\tilde{p}}$ we get:

$$\|y - \tilde{X}b\|_2^2 = \|\tilde{X}_C^\top (y - \tilde{X}b)\|_2^2 = \left\| \begin{bmatrix} \tilde{X}^\top \\ C^\top \end{bmatrix} y - \begin{bmatrix} b \\ \mathbf{0} \end{bmatrix} \right\|_2^2 = \|\tilde{X}^\top y - b\|_2^2 + \text{const}, \quad (\text{C.1})$$

which implies that under orthogonal situation the optimization problem in (2.5) could be recast as

$$\arg \min_b \left\{ \frac{1}{2} \|\tilde{y} - b\|_2^2 + \sigma J_\lambda(W\llbracket b \rrbracket_{\mathbb{I}}) \right\}, \quad (\text{C.2})$$

for $\tilde{y} := \tilde{X}^\top y$. After introducing new variable to problem (C.2), namely $c \in \mathbb{R}^m$, we get the equivalent formulation

$$\arg \min_{b,c} \left\{ \frac{1}{2} \|\tilde{y} - b\|_2^2 + \sigma J_\lambda(c) : c = W\llbracket b \rrbracket_{\mathbb{I}} \right\}. \quad (\text{C.3})$$

Proposition C.1. *Let $f(b, c) : \mathbb{R}^p \times \mathbb{R}^m \rightarrow \mathbb{R}$ be any function and consider optimization problem $\arg \min_{b,c} \{f(b, c) : (b, c) \in \mathcal{D}\}$ with unique solution (b^*, c^*) and feasible set $\mathcal{D} \subset \mathbb{R}^p \times \mathbb{R}^m$.*

Define $\mathcal{D}^c := \{c \in \mathbb{R}^m \mid \exists b \in \mathbb{R}^p : (b, c) \in \mathcal{D}\}$. Suppose that for any $c \in \mathcal{D}^c$, there exists unique solution, b^c , to problem $\arg \min_b \{f(b, c) : (b, c) \in \mathcal{D}\}$. Moreover, assume that the solution to $\arg \min_c \{f(b^c, c) : c \in \mathcal{D}^c\}$ is unique. Then, it occurs

$$\begin{cases} c^* = \arg \min_c \{f(b^c, c) : c \in \mathcal{D}^c\} \\ b^* = b^{c^*} \end{cases}. \quad (\text{C.4})$$

Proof. Suppose that there exists $(b^0, c^0) \in \mathcal{D}$, such that $f(b^0, c^0) < f(b^*, c^*)$, where b^* and c^* are defined as in (C.4). We have

$$f(b^{c^0}, c^0) \leq f(b^0, c^0) < f(b^*, c^*) = f(b^{c^*}, c^*), \quad (\text{C.5})$$

which leads to the contradiction with definition of c^* . \blacksquare

We will apply the above proposition to (C.3). Let (b^*, c^*) be solution to (C.3). Then b^* is also solution to convex problem (C.2) with strictly convex objective function and therefore is unique. Since $c^* = W\llbracket b^* \rrbracket_{\mathbb{I}}$, c^* is unique as well. In considered situation $\mathcal{D}^c = \{c : c \succeq 0\}$. We will start with solving the problem $b^c = \arg \min_b \left\{ \frac{1}{2} \|\tilde{y} - b\|_2^2 + \sigma J_\lambda(c) : c = W\llbracket b \rrbracket_{\mathbb{I}} \right\}$. The additive constant in the objective could be omitted. Moreover, for each $i \in \{1, \dots, m\}$ we have

$$b_{\mathbb{I}_i}^c = \arg \min_{b_{\mathbb{I}_i}} \left\{ \|\tilde{y}_{\mathbb{I}_i} - b_{\mathbb{I}_i}\|_2^2 : w_i^2 \|b_{\mathbb{I}_i}\|_2^2 - c_i^2 = 0 \right\}. \quad (\text{C.6})$$

The Lagrange Multipliers method quickly yields $b_{\mathbb{I}_i}^c = (w_i \|\tilde{y}_{\mathbb{I}_i}\|_2)^{-1} c_i \tilde{y}_{\mathbb{I}_i}$ and, consequently, it holds $\|\tilde{y}_{\mathbb{I}_i} - b_{\mathbb{I}_i}^c\|_2^2 = (\|\tilde{y}_{\mathbb{I}_i}\|_2 - w_i^{-1} c_i)^2$. From Proposition C.1, we get the following procedure for solution, b^* , to problem (C.2)

$$\begin{cases} c^* = \arg \min_c \left\{ \frac{1}{2} \sum_{i=1}^m (\|\tilde{y}_{\mathbb{I}_i}\|_2 - w_i^{-1} c_i)^2 + J_{\sigma\lambda}(c) \right\} \\ b_{\mathbb{I}_i}^* = c_i^* (w_i \|\tilde{y}_{\mathbb{I}_i}\|_2)^{-1} \tilde{y}_{\mathbb{I}_i}, \quad i = 1, \dots, m \end{cases} \quad (\text{C.7})$$

(notice that we applied Proposition D.2 to omit the constraints $c \succeq 0$ and that the objective function in definition of c^* is strictly feasible, which guarantees the unique solution. The above procedure yields conclusion, that indices of groups estimated by gSLOPE as relevant coincide with the support of solution to SLOPE problem with diagonal matrix having inverses of weights w_1, \dots, w_m on diagonal. Moreover, after defining $\tilde{\beta} \in \mathbb{R}^{\tilde{p}}$ by conditions $\tilde{\beta}_{\mathbb{I}_i} := R_i \beta_{I_i}$, $i = 1, \dots, m$, we simply have $\llbracket \tilde{\beta} \rrbracket_{\mathbb{I}} = \llbracket \beta \rrbracket_{X, I}$ and

$$\tilde{y} = \tilde{X}^\top y = \tilde{X}^\top \left(\sum_{i=1}^m U_i R_i \beta_{I_i} + z \right) = \tilde{X}^\top (\tilde{X} \tilde{\beta} + z) = \tilde{\beta} + \tilde{X}^\top z, \quad \text{hence } \tilde{y} \sim \mathcal{N}(\tilde{\beta}, \sigma^2 \mathbf{I}_{\tilde{p}}). \quad (\text{C.8})$$

Summarizing, if the assumption about the orthogonality at groups level is in use, one can consider the statistically equivalent model $\tilde{y} \sim \mathcal{N}(\tilde{\beta}, \sigma^2 \mathbf{I}_{\tilde{p}})$, define truly relevant groups via the support of $\llbracket \tilde{\beta} \rrbracket_{\mathbb{I}}$ and treat the vector $\llbracket b^* \rrbracket_{\mathbb{I}} = (\frac{c_1^*}{w_1}, \dots, \frac{c_m^*}{w_m})$ as an gSLOPE estimate of group effect sizes, where b^* and c^* are defined in (2.8), i.e. it holds $\llbracket b^* \rrbracket_{\mathbb{I}} = \llbracket \beta^{\text{gs}} \rrbracket_{X, I}$ for any solution β^{gs} to problem (2.2).

D SLOPE with diagonal experiment matrix

Let $y \in \mathbb{R}^p$ be fixed vector and d_1, \dots, d_p be positive numbers. We will use notation $\text{diag}(d_1, \dots, d_p)$ to define the diagonal matrix D such as $D_{i,i} = d_i$ for $i = 1, \dots, p$. Denote $d := (d_1, \dots, d_p)^\top$ and let b^* be the solution to SLOPE optimization problem with diagonal experiment matrix, i.e. the solution to

$$\underset{b}{\text{minimize}} \quad f(b) := \frac{1}{2} \|y - Db\|_2^2 + J_\lambda(b). \quad (\text{D.1})$$

Since f is strictly convex function, the solution to (D.1) is unique. It is easy to observe, that changing sign of y_i corresponds to changing sign at i th coefficient of solution as well as permuting coefficients of y together with d_i 's permutes coefficients of b^* . We will summarize this observations below without proofs.

Proposition D.1. *Let $\pi : \{1, \dots, p\} \rightarrow \{1, \dots, p\}$ be given permutation with P_π as corresponding matrix. Then:*

i) $P_\pi D P_\pi^\top = \text{diag}(d_{\pi(1)}, \dots, d_{\pi(p)})$;

ii) $b_\pi := P_\pi b^*$ is solution to minimize $f_\pi(b) := \frac{1}{2} \|P_\pi y - P_\pi D P_\pi^\top b\|_2^2 + J_\lambda(b)$;

iii) $b_S := S b^*$ is solution to minimize $f_S(b) := \frac{1}{2} \|S y - D b\|_2^2 + J_\lambda(b)$,

where S is diagonal matrix with entries on diagonal coming from set $\{-1, 1\}$.

Proposition D.2. *If $y \succeq 0$, then $b^* \succeq 0$.*

Proof. Suppose that for some r it occurs $b_r < 0$ for any $b \in \mathbb{R}^p$. If $y_r = 0$, then taking \hat{b} defined as $\hat{b}_i := \begin{cases} 0, & i = r \\ b_i, & \text{otherwise} \end{cases}$, we get $|\hat{b}| \leq |b|$ and Corollary A.2 gives $J_\lambda(\hat{b}) \leq J_\lambda(b)$. Consequently,

$$f(b) - f(\hat{b}) \geq \frac{1}{2} \|y - Db\|_2^2 - \frac{1}{2} \|y - D\hat{b}\|_2^2 = \frac{1}{2} (y_r - d_r b_r)^2 - \frac{1}{2} (y_r + d_r \hat{b}_r)^2 = \frac{1}{2} d_r^2 b_r^2 > 0.$$

Hence b could not be the solution. Now consider case when $y_r > 0$ and define \hat{b} by putting

$$\hat{b}_i := \begin{cases} -b_r, & i = r \\ b_i, & \text{otherwise} \end{cases}. \text{ Then we have } J_\lambda(b) = J_\lambda(\hat{b}) \text{ and}$$

$$f(b) - f(\hat{b}) = \frac{1}{2} (y_r - d_r b_r)^2 - \frac{1}{2} (y_r + d_r b_r)^2 = -2y_r d_r b_r > 0.$$

and, as before, b could not be optimal. ■

Proposition D.3. *Let b^* be the solution to problem (D.1), $\{y_i\}_{i=1}^p$ be nonnegative sequence, $\{d_i\}_{i=1}^p$ be the sequence of positive numbers and assume that*

$$d_1 y_1 \geq \dots \geq d_p y_p. \quad (\text{D.2})$$

If b^ has exactly r nonzero entries for $r > 0$, then the set $\{1, \dots, r\}$ corresponds to the support of b^* .*

Proof. It is enough to show that

$$(j \in \{2, \dots, m\}, b_j^* \neq 0) \implies b_{j-1}^* \neq 0.$$

Suppose that this is not true. From Proposition D.2 we know that b^* is nonnegative, hence we can find i from $\{2, \dots, m\}$ such as $b_j^* > 0$ and $b_{j-1}^* = 0$. For $\varepsilon \in (0, b_j^*/2]$ define vector b_ε by putting $(b_\varepsilon)_{j-1} := \varepsilon$, $(b_\varepsilon)_j := b_j^* - \varepsilon$ and $(b_\varepsilon)_i := b_i^*$ for $i \notin \{j, l\}$. From Proposition A.3 we have that $J_\lambda(b_\varepsilon) \leq J_\lambda(b^*)$, which gives

$$\begin{aligned} f(b^*) - f(b_\varepsilon) &\geq \frac{1}{2}(y_{j-1} - d_{j-1}b_{j-1}^*)^2 + \frac{1}{2}(y_j - d_j b_j^*)^2 \\ &\quad - \frac{1}{2}(y_{j-1} - d_{j-1}b_\varepsilon(j-1))^2 - \frac{1}{2}(y_j - d_j b_\varepsilon(j))^2 = \\ &\quad \varepsilon \left(A - \frac{d_{j-1}^2 + d_j^2}{2} \cdot \varepsilon \right), \end{aligned} \tag{D.3}$$

for $A := (y_{j-1}d_{j-1} - y_j d_j) + d_j^2 b_j^* > 0$.

Therefore, we can find $\varepsilon > 0$ such as $f(b^*) > f(b_\varepsilon)$, which contradicts the optimality of b^* . \blacksquare

Consider now problem (D.1) with arbitrary sequence $\{y_i\}_{i=1}^p$. Suppose that b^* has exactly $r > 0$ nonzero coefficients and that $\pi : \{1, \dots, p\} \rightarrow \{1, \dots, p\}$ is permutation which gives the order of magnitudes for Dy , i.e. $d_{\pi(1)}|y|_{\pi(1)} \geq \dots \geq d_{\pi(p)}|y|_{\pi(p)}$. Basing on our previous observations, we get important

Corollary D.4. *If b^* is the solution to (D.1) having exactly $r > 0$ nonzero coefficients and π is permutation which places components of $D|y|$ in a nonincreasing order, i.e. $d_{\pi(i)}|y|_{\pi(i)} = |Dy|_{(i)}$ for $i = 1, \dots, p$, then the support of b^* is composed of the set $\{\pi(1), \dots, \pi(r)\}$.*

The next three lemmas were proven in [10] in situation when $d_1 = \dots = d_p = 1$. We will follow the reasoning from this paper to prove the generalized claims. The main difference is that in general case the solution to considered problem (D.1) does not have to be nonincreasingly ordered, under assumption that $d_1 y_1 \geq \dots \geq d_p y_p \geq 0$ (which is the case for $d_1 = \dots = d_p = 1$). This makes that generalizations of proofs presented in [10] are not straightforward.

Lemma D.5. *Consider nonnegative sequence $\{y_i\}_{i=1}^p$ and sequence of positive numbers $\{d_i\}_{i=1}^p$ such as $d_1 y_1 \geq \dots \geq d_p y_p$. If b^* is solution to problem (D.1) having exactly r nonzero entries, then for every $j \leq r$ it holds that*

$$\sum_{i=j}^r (d_i y_i - \lambda_i) > 0 \tag{D.4}$$

and for every $j \geq r+1$

$$\sum_{i=r+1}^j (d_i y_i - \lambda_i) \leq 0. \tag{D.5}$$

Proof. From Proposition D.3 we know that $b_i^* > 0$ for $i \in \{1, \dots, r\}$. Let us define

$$\tilde{b}_i := \begin{cases} b_i^* - h, & i \in \{j, \dots, r\} \\ b_i^*, & \text{otherwise.} \end{cases},$$

where we restrict only to sufficiently small values of h , so as to the condition $\tilde{b}_i > 0$ is met for all i from $\{j, \dots, r\}$. For such h we have $b_{(r+1)}^* = \dots = b_{(p)}^* = \tilde{b}_{(r+1)} = \dots = \tilde{b}_{(p)} = 0$. Therefore

there exists permutation $\pi : \{1, \dots, r\} \rightarrow \{1, \dots, r\}$ such as $\sum_{i=1}^r \lambda_i \tilde{b}_{(i)} = \sum_{i=1}^r \lambda_{\pi(i)} \tilde{b}_i$. For such permutation we have

$$\begin{aligned} J_\lambda(b^*) - J_\lambda(\tilde{b}) &= \sum_{i=1}^r \lambda_i b_{(i)}^* - \sum_{i=1}^r \lambda_i \tilde{b}_{(i)} = \sum_{i=1}^r \lambda_i b_{(i)}^* - \sum_{i=1}^r \lambda_{\pi(i)} \tilde{b}_i \\ &\geq \sum_{i=1}^r \lambda_{\pi(i)} b_i^* - \sum_{i=1}^r \lambda_{\pi(i)} \tilde{b}_i = h \sum_{i=j}^r \lambda_{\pi(i)} \geq h \sum_{i=j}^r \lambda_i, \end{aligned} \quad (\text{D.6})$$

where the first inequality follows from the rearrangement inequality and second is the consequence of monotonicity of $\{\lambda_i\}_{i=1}^p$. We also have

$$\begin{aligned} \|y - Db^*\|_2^2 - \|y - D\tilde{b}\|_2^2 &= \sum_{i=j}^r (y_i - d_i b_i^*)^2 - \sum_{i=j}^r (y_i - d_i b_i^* + d_i h)^2 \\ &= 2h \sum_{i=j}^r (d_i^2 b_i^* - d_i y_i) - h^2 \sum_{i=j}^r d_i^2. \end{aligned} \quad (\text{D.7})$$

Optimality of b^* , (D.6) and (D.7) yield

$$0 \geq f(b^*) - f(\tilde{b}) \geq h \sum_{i=j}^r (d_i^2 b_i^* - d_i y_i + \lambda_i) - \frac{1}{2} h^2 \sum_{i=j}^r d_i^2, \quad (\text{D.8})$$

for each h from the interval $[0, \varepsilon]$, where $\varepsilon > 0$ is some (sufficiently small) value. This gives $\sum_{i=j}^r (d_i^2 b_i^* - d_i y_i + \lambda_i) \leq 0$ and consequently

$$\sum_{i=j}^r (d_i y_i - \lambda_i) \geq \sum_{i=j}^r d_i^2 b_i^* > 0. \quad (\text{D.9})$$

To prove claim (D.5), consider a new sequence defined as $\tilde{b}_i := \begin{cases} h, & i \in \{r+1, \dots, j\} \\ b_i^*, & \text{otherwise.} \end{cases}$. We will

restrict our attention only to $0 < h < \min\{b_i^* : i \leq r\}$, so as to $b_{(\cdot)}^*$ and $\tilde{b}_{(\cdot)}$ are given by applying the same permutation to b^* and \tilde{b} , respectively. Moreover, for each i from $\{r+1, \dots, j\}$ it holds $\tilde{b}_{(i)} = \tilde{b}_i = h$. From optimality of b^*

$$0 \geq f(b^*) - f(\tilde{b}) = \frac{1}{2} \sum_{i=r+1}^j (y_i^2 - (y_i - d_i h)^2) - \sum_{i=r+1}^j \lambda_i h = h \sum_{i=r+1}^j (d_i y_i - \lambda_i) - \frac{1}{2} h^2 \sum_{i=r+1}^j d_i^2,$$

for all considered h , which leads to (D.5). ■

Lemma D.6. *Let b^* be solution to problem (D.1) with nonnegative, nonincreasing sequence $\{\lambda_i\}_{i=1}^p$. Let $R(b^*)$ be number of all nonzeros in b^* and $r \geq 1$. Then, for any $i \in \{1, \dots, p\}$*

$$\{y : b_i^* \neq 0 \text{ and } R(b^*) = r\} = \{y : d_i |y_i| > \lambda_r \text{ and } R(b^*) = r\}.$$

Proof. Suppose that b^* has $r > 0$ nonzero coefficients and let π be permutation which places components of $D|y|$ in a nonincreasing order. From Corollary D.4 it holds that $\{i : b_i^* \neq 0\} = \{\pi(1), \dots, \pi(r)\}$. Define $\tilde{y} := P_\pi S y$ and $\tilde{D} := P_\pi D P_\pi^\top$, for S being the diagonal matrix such as $S_{i,i} = \text{sgn}(y_i)$. Then $P_\pi S b^*$ is solution to problem

$$\arg \min_b \frac{1}{2} \|\tilde{y} - \tilde{D}b\|_2^2 + J_\lambda(b), \quad (\text{D.10})$$

which satisfies the assumptions of Lemma D.5. Taking $j = r$ in (D.4) and $j = r + 1$ in (D.5) we immediately get

$$d_{\pi(r)}|y|_{\pi(r)} > \lambda_r \quad \text{and} \quad d_{\pi(r+1)}|y|_{\pi(r+1)} \leq \lambda_{r+1}. \quad (\text{D.11})$$

We will now show that $\{y : b_i^* \neq 0 \text{ and } R(b^*) = r\} \subset \{y : d_i|y_i| > \lambda_r \text{ and } R(b^*) = r\}$. Fix $i \in \{1, \dots, p\}$ and suppose that b_i^* is nonzero coefficient. Then $i \in \{\pi(1), \dots, \pi(r)\}$ and therefore $d_i|y_i| \geq d_{\pi(r)}|y|_{\pi(r)} > \lambda_r$, thanks to first inequality from (D.11). To show the second inclusion assume that $d_i|y_i| > \lambda_r$. Then, from the second inequality in (D.11), $d_i|y_i| > \lambda_{r+1} \geq d_{\pi(r+1)}|y|_{\pi(r+1)}$, which gives $i \in \{\pi(1), \dots, \pi(r)\}$. \blacksquare

Lemma D.7. *For given sequence $\{y_i\}_{i=1}^p$, sequence of positive numbers $\{d_i\}_{i=1}^p$, nonincreasing, nonnegative sequence $\{\lambda_i\}_{i=1}^p$ and fixed $j \in \{1, \dots, p\}$, consider a following procedure*

- define $\tilde{y} := (y_1, \dots, y_{j-1}, y_{j+1}, \dots, y_p)^\top$, $\tilde{D} := \text{diag}(d_1, \dots, d_{j-1}, d_{j+1}, \dots, d_p)$, $\tilde{d}_i := \tilde{D}_{i,i}$ for $i = 1, \dots, p-1$ and $\tilde{\lambda} := (\lambda_2, \dots, \lambda_p)^\top$;
- find $\tilde{b}^* := \arg \min_{b \in \mathbb{R}^{p-1}} \frac{1}{2} \|\tilde{y} - \tilde{D}b\|_2^2 + J_{\tilde{\lambda}}(b)$;
- define $\tilde{R}^j(\tilde{b}^*) := |\{i : \tilde{b}_i^* \neq 0\}|$.

Then for $r \geq 1$ it holds $\{y : d_j|y_j| > \lambda_r \text{ and } R(b^*) = r\} \subset \{y : d_j|y_j| > \lambda_r \text{ and } \tilde{R}^j(\tilde{b}^*) = r-1\}$.

Proof. We have to show that solution \tilde{b}^* to problem

$$\underset{b}{\text{minimize}} \quad F(b) := \frac{1}{2} \sum_{i=1}^{p-1} (\tilde{y}_i - \tilde{d}_i b_i)^2 + \sum_{i=1}^{p-1} \tilde{\lambda}_i b_{(i)} \quad (\text{D.12})$$

has exactly $r-1$ nonzero coefficients. From Proposition D.1 we know that the change of signs of y_i 's does not affect the support, hence without loss of generality we can assume that $\tilde{y} \geq 0$, and $\tilde{b}^* \geq 0$ as a result (from Proposition D.2). We will start with situation when $d_1 y_1 \geq \dots \geq d_p y_p$ and consequently $\tilde{d}_1 \tilde{y}_1 \geq \dots \geq \tilde{d}_{p-1} \tilde{y}_{p-1}$. If j is fixed index such as $d_j|y_j| > \lambda_r$ and $R(b^*) = r$, this gives

$$j \in \{1, \dots, r\}. \quad (\text{D.13})$$

To show that solution to (D.12) has at least $r-1$ nonzero entries, suppose by contradiction that \tilde{b}^* has exactly $k-1$ nonzero entries with $k < r$. Let us define $\hat{b} \in \mathbb{R}^{p-1}$ as

$$\hat{b}_i := \begin{cases} h, & i \in \{k, \dots, r-1\} \\ \tilde{b}_i^*, & \text{otherwise} \end{cases},$$

where $0 < h < \min\{\tilde{b}_1^*, \dots, \tilde{b}_{k-1}^*\}$. Then

$$F(\tilde{b}^*) - F(\hat{b}) = h \sum_{i=k}^{r-1} (\tilde{d}_i \tilde{y}_i - \tilde{\lambda}_i) - h^2 \sum_{i=k}^{r-1} \frac{1}{2} \tilde{d}_i^2. \quad (\text{D.14})$$

Now

$$\sum_{i=k}^{r-1} (\tilde{d}_i \tilde{y}_i - \tilde{\lambda}_i) = \sum_{i=k+1}^r (\tilde{d}_{i-1} \tilde{y}_{i-1} - \lambda_i) \geq \sum_{i=k+1}^r (d_i y_i - \lambda_i) > 0, \quad (\text{D.15})$$

where the first equality follows from $\tilde{\lambda}_i = \lambda_{i+1}$, the first inequality from $\tilde{d}_{i-1} \tilde{y}_{i-1} \geq d_i y_i$ and the second from Lemma D.5. If h is small enough, we get $F(\hat{b}) < F(\tilde{b}^*)$ which leads to contradiction.

Suppose now by contradiction that \tilde{b}^* has k nonzero entries with $k \geq r$ and define

$$\hat{b}_i := \begin{cases} \tilde{b}_i^* - h, & i \in \{r, \dots, k\} \\ \tilde{b}_i^*, & \text{otherwise} \end{cases}.$$

Analogously to (D.6), we get $J_{\tilde{\lambda}}(\tilde{b}^*) - J_{\tilde{\lambda}}(\hat{b}) \geq h \sum_{i=r}^k \tilde{\lambda}_i$ and consequently

$$F(\tilde{b}^*) - F(\hat{b}) \geq h \left[\sum_{i=r}^k (\tilde{\lambda}_i - \tilde{d}_i \tilde{y}_i) + \sum_{i=r}^k \tilde{d}_i^2 \tilde{b}_i^* \right] - \frac{1}{2} h^2 \sum_{i=r}^k \tilde{d}_i^2. \quad (\text{D.16})$$

Now

$$\sum_{i=r}^k (\tilde{\lambda}_i - \tilde{d}_i \tilde{y}_i) = \sum_{i=r+1}^{k+1} (\lambda_i - d_i y_i) \geq 0, \quad (\text{D.17})$$

where the first equality follows from definition of $\tilde{\lambda}$ and (D.13), while the inequality follows from Lemma D.5. If h is small enough, we get $F(\hat{b}) < F(\tilde{b}^*)$, which contradicts the optimality of \tilde{b}^* .

Consider now general situation, i.e. without assumption concerning the order of $D|y|$. Suppose that π , with corresponding matrix P_π , is permutation which orders $D|y|$. Define $y_\pi := P_\pi y$ and $D_\pi := P_\pi D P_\pi^\top$. Applying the procedure described in the statement of Lemma simultaneously to (y, D, λ) for j , and to (y_π, D_π, λ) for $\pi(j)$ we end with $(\tilde{y}, \tilde{D}, \tilde{\lambda}, \tilde{R}_1^j(\tilde{b}^*))$ and $(\tilde{y}_\pi, \tilde{D}_\pi, \tilde{\lambda}, \tilde{R}_2^{\pi(j)}(\tilde{b}_\pi^*))$. It is straightforward to see, that there exists permutation $\tilde{\pi} : \{1, \dots, p-1\} \rightarrow \{1, \dots, p-1\}$ such that $\tilde{y}_\pi = P_{\tilde{\pi}} \tilde{y}$ and $\tilde{D}_\pi = P_{\tilde{\pi}} \tilde{D} P_{\tilde{\pi}}^\top$. From Proposition D.1 we have that $\tilde{b}_\pi^* = P_{\tilde{\pi}} \tilde{b}^*$ and $\tilde{R}_1^j(\tilde{b}^*) = \tilde{R}_2^{\pi(j)}(\tilde{b}_\pi^*)$. Moreover, from the first part of proof $\tilde{R}_2^{\pi(j)}(\tilde{b}_\pi^*) = r-1$, which gives the claim. \blacksquare

E Minimax estimation of gSLOPE

Proof of Theorem 2.6. Once again we will employ the equivalent formulation of gSLOPE under assumption about orthogonality at groups level, i.e. problem (2.8), and we will consider statistically equivalent model $\tilde{y} \sim \mathcal{N}(\tilde{\beta}, \sigma^2 \mathbf{I}_{\tilde{p}})$, with $\tilde{\beta}_{\mathbb{I}_i} = R_i \beta_{I_i}$, $i = 1, \dots, m$. Then $\llbracket \beta \rrbracket_{X,I} = \llbracket \tilde{\beta} \rrbracket_{\mathbb{I}}$ and for solution b^* to (2.8) it holds $\llbracket b^* \rrbracket_{\mathbb{I}} = \llbracket \beta^{\text{gs}} \rrbracket_{X,I}$ for any solution β^{gs} to problem (2.2). Without loss of generality, assume $\sigma = 1$. Note that $\|\tilde{y}_{\mathbb{I}_i}\|_2^2$ is distributed as the noncentral $\chi_{l_i}^2(\|\tilde{\beta}_{\mathbb{I}_i}\|_2^2)$, where $\|\tilde{\beta}_{\mathbb{I}_i}\|_2^2$ is the noncentrality.

The lower bound of the minimax risk can be obtained as follows. For each \mathbb{I}_i , only $\tilde{\beta}_j$ with the smallest index $j \in \mathbb{I}_i$ is *possibly* nonzero and the rest $l_i - 1$ components of $\tilde{\beta}_{\mathbb{I}_i}$ are fixed to be zero. Then, this is reduced to a simple Gaussian sequence model with length m and sparsity at most k . Given the condition $k/m \rightarrow 0$, this classical sequence model has minimax risk $(1 + o(1))2k \log(m/k)$ (see e.g. [13]).

Our next step is to evaluate the worst risk of gSLOPE over the nearly black object. We would complete the proof if we show this worst risk is bounded above by $(1 + o(1))2k \log(m/k)$. For simplicity, assume that $\|\tilde{\beta}_{\mathbb{I}_i}\|_2 = 0$ for all $i \geq k+1$ and write $\mu_i = \|\tilde{\beta}_{\mathbb{I}_i}\|_2$, $\zeta_i = \|\tilde{y}_{\mathbb{I}_i}\|_2 \sim \chi_{l_i}(\mu_i)$. Denote by $\hat{\zeta}$ the SLOPE solution. Then, the risk is

$$\mathbb{E} \|\hat{\zeta} - \mu\|_2^2 = \mathbb{E} \sum_{i=1}^k (\hat{\zeta}_i - \mu_i)^2 + \mathbb{E} \sum_{i=k+1}^m \hat{\zeta}_i^2.$$

Then, it suffices to show

$$\mathbb{E} \left[\sum_{i=1}^k (\hat{\zeta}_i - \mu_i)^2 \right] \leq (1 + o(1))2k \log(m/k) \quad (\text{E.1})$$

and

$$\mathbb{E} \left[\sum_{i=k+1}^m \widehat{\zeta}_i^2 \right] = o(1)2k \log(m/k). \quad (\text{E.2})$$

Below, Lemmas E.1, E.2, and E.3 together give (E.2). The remaining part of this proof serves to validate (E.1). To start with, we employ the representation $\zeta_i^2 = (\xi_{i1} + \mu_i)^2 + \xi_{i2}^2 + \dots + \xi_{il_i}^2$ for i.i.d. $\xi_{ij} \sim \mathcal{N}(0, 1)$ (we can assume this representation without loss of generality, since the distribution of $(\xi_{i1} + a_1)^2 + (\xi_{i2} + a_2)^2 + \dots + (\xi_{il_i} + a_{l_i})^2$ depends only on the non-centrality $a_1^2 + \dots + a_{l_i}^2$). As in the proof of Lemma 3.2 in [26], we get

$$\begin{aligned} \sum_{i=1}^k (\widehat{\zeta}_i - \mu_i)^2 &\leq \left(\|\widehat{\zeta}_{[1:k]} - \zeta_{[1:k]}\|_2 + \|\zeta_{[1:k]} - \mu_{[1:k]}\|_2 \right)^2 \\ &\leq \left(\|\lambda_{[1:k]}\|_2 + \|\zeta_{[1:k]} - \mu_{[1:k]}\|_2 \right)^2. \end{aligned} \quad (\text{E.3})$$

As l is fixed and $k/m \rightarrow 0$, [18] gives $\lambda_i \sim \sqrt{2 \log \frac{m}{qi}}$ for all $i \leq k$. From this we know

$$\|\lambda_{[1:k]}\|_2^2 = \sum_{i=1}^k \lambda_i^2 \sim 2k \log \frac{m}{k}. \quad (\text{E.4})$$

Next, we see

$$\begin{aligned} \left| \sqrt{(\xi_{i1} + \mu_i)^2 + \xi_{i2}^2 + \dots + \xi_{il_i}^2} - \mu_i \right| &\leq \sqrt{\xi_{i2}^2 + \dots + \xi_{il_i}^2} + |\xi_{i1}| \\ &\leq 2\sqrt{\xi_{i1}^2 + \xi_{i2}^2 + \dots + \xi_{il_i}^2} \equiv 2\|\xi_i\|_2, \end{aligned}$$

which yields

$$\|\zeta_{[1:k]} - \mu_{[1:k]}\|_2^2 \leq 4 \sum_{i=1}^k \|\xi_i\|_2^2 \quad (\text{E.5})$$

Note that $\sum_{i=1}^k \|\xi_i\|_2^2$ is distributed as the chi-square with $l_1 + \dots + l_k \leq lk$ degrees of freedom. Taking (E.4) and (E.5) together, from (E.3) we get

$$\begin{aligned} \mathbb{E} \left[\sum_{i=1}^k (\widehat{\zeta}_i - \mu_i)^2 \right] &\leq \|\lambda_{[1:k]}\|_2^2 + \mathbb{E} \|\zeta_{[1:k]} - \mu_{[1:k]}\|_2^2 + 2\|\lambda_{[1:k]}\|_2 \mathbb{E} \|\zeta_{[1:k]} - \mu_{[1:k]}\|_2 \\ &\leq (1 + o(1))2k \log \frac{m}{k} + 4lk + 2\sqrt{(1 + o(1))2k \log \frac{m}{k}} \cdot \sqrt{4lk} \\ &\sim (1 + o(1))2k \log \frac{m}{k}, \end{aligned}$$

where the last step makes use of $m/k \rightarrow \infty$. This establishes (E.1) and consequently completes the proof. ■

The following three lemmas aim to prove (E.2). Denote by $\zeta_{(1)} \geq \dots \geq \zeta_{(m-k)}$ the order statistics of $\zeta_{k+1}, \dots, \zeta_m$. Recall that $\zeta_i \sim \chi_{l_i}$ for $i \geq k+1$. As in the proof of Lemma 3.3 in [26], we have

$$\sum_{i=k+1}^m \widehat{\zeta}_i^2 \leq \sum_{i=1}^{m-k} (\zeta_{(i)} - \lambda_{k+i})_+^2,$$

where $x_+ = \max\{x, 0\}$. For a sufficiently large constant $A > 0$ and sufficiently small constant $\alpha > 0$ both to be specified later, we partition the sum into three parts:

$$\sum_{i=1}^{m-k} \mathbb{E}(\zeta_{(i)} - \lambda_{k+i})_+^2 = \sum_{i=1}^{\lfloor Ak \rfloor} \mathbb{E}(\zeta_{(i)} - \lambda_{k+i})_+^2 + \sum_{i=\lfloor Ak \rfloor}^{\lfloor \alpha m \rfloor} \mathbb{E}(\zeta_{(i)} - \lambda_{k+i})_+^2 + \sum_{i=\lfloor \alpha m \rfloor}^{m-k} \mathbb{E}(\zeta_{(i)} - \lambda_{k+i})_+^2$$

The three lemmas, respectively, show that each part is negligible compared with $2k \log(m/k)$. We indeed prove a stronger version in which the order statistics $\zeta_{(1)} \geq \dots \geq \zeta_{(m-k)} \geq \zeta_{(m-k+1)} \geq \dots \geq \zeta_{(m)}$ come from m i.i.d. χ_l . Let U_1, \dots, U_m be i.i.d. uniform random variables on $(0, 1)$, and $U_{(1)} \leq U_{(2)} \leq \dots \leq U_{(m)}$ be the *increasing* order statistics. So we have the representation $\zeta_{(i)} = F_{\chi_l}^{-1}(1 - U_{(i)})$

Lemma E.1. *Under the preceding conditions, for any $A > 0$ we have*

$$\frac{1}{2k \log(m/k)} \sum_{i=1}^{\lfloor Ak \rfloor} \mathbb{E}(\zeta_{(i)} - \lambda_{k+i})_+^2 \rightarrow 0.$$

Proof of Lemma E.1. Recognizing that l is fixed, from [18] it follows that

$$F_{\chi_l}^{-1}(1 - q_1) - F_{\chi_l}^{-1}(1 - q_2) \sim \sqrt{2 \log \frac{1}{q_1}} - \sqrt{2 \log \frac{1}{q_2}}$$

for $q_1, q_2 \rightarrow 0$. We also know that ζ_i is distributed as $F_{\chi_l}^{-1}(1 - U_{(i)})$. Making use of these facts, we get

$$\begin{aligned} \mathbb{E}(\zeta_{(i)} - \lambda_{k+i})_+^2 &= \mathbb{E}(F_{\chi_l}^{-1}(1 - U_{(i)}) - F_{\chi_l}^{-1}(1 - q(k+i)/m))^2 \\ &\sim \mathbb{E} \left(\sqrt{2 \log \frac{1}{U_{(i)}}} - \sqrt{2 \log \frac{m}{q(k+i)}} \right)_+^2 \\ &\leq \mathbb{E} \left(\sqrt{2 \log \frac{1}{U_{(i)}}} - \sqrt{2 \log \frac{m}{q(k+i)}} \right)^2 \\ &\lesssim \mathbb{E} \left(\frac{\log^2(q(k+i)/mU_{(i)})}{\log(m/q(k+i))} \right). \end{aligned}$$

Now, we proceed to evaluate

$$\mathbb{E} \left[\log^2 \frac{q(k+i)}{mU_{(i)}} \right] = \log^2 \frac{q(k+i)}{m} + \mathbb{E} \log^2 U_{(i)} - 2 \log \frac{q(k+i)}{m} \mathbb{E} \log U_{(i)}.$$

Observing that $U_{(i)}$ follows Beta($i, m+i-i$), we get (see e.g. [2])

$$\begin{aligned} \mathbb{E} \log U_{(i)} &= -\log \frac{m+1}{i} + \delta_1, \\ \mathbb{E} \log^2 U_{(i)} &= \left(\log \frac{m+1}{i} - \delta_1 \right)^2 + \frac{1}{i} - \frac{1}{m+1} + \delta_2 \end{aligned}$$

for some $\delta_1 = O(1/i)$ and $\delta_2 = O(1/i^2)$. Thus we can evaluate $\mathbb{E} \log^2 \frac{q(k+i)}{mU_{(i)}}$ as

$$\begin{aligned} \mathbb{E} \log^2 \frac{q(k+i)}{mU_{(i)}} &= \log^2 \frac{q(k+i)}{m} - 2 \log \frac{q(k+i)}{m} \mathbb{E} \log U_{(i)} + \mathbb{E} \log^2 U_{(i)} \\ &= \log^2 \frac{q(k+i)}{m} + 2 \log \frac{q(k+i)}{m} \left(\log \frac{m+1}{i} - \delta_1 \right) + \left(\log \frac{m+1}{i} - \delta_1 \right)^2 + \frac{1}{i} - \frac{1}{m+1} + \delta_2 \\ &= \log^2 \frac{q(k+i)(m+1)}{im} - 2\delta_1 \log \frac{q(k+i)(m+1)}{im} + \frac{1}{i} - \frac{1}{m+1} + \delta_1^2 + \delta_2. \end{aligned}$$

Hence, we get

$$\begin{aligned}
& \sum_{i=1}^{\lfloor Ak \rfloor} \mathbb{E}(\zeta_{(i)} - \lambda_{k+i})_+^2 \\
& \lesssim \frac{1}{\log \frac{m}{q(A+1)k}} \left(\underbrace{\sum_{i=1}^{\lfloor Ak \rfloor} \log^2 \frac{q(k+i)(m+1)}{im}}_{\text{I}} - \underbrace{\sum_{i=1}^{\lfloor Ak \rfloor} 2\delta_1 \log \frac{q(k+i)(m+1)}{im}}_{\text{II}} + \underbrace{\sum_{i=1}^{\lfloor Ak \rfloor} \left(\frac{1}{i} - \frac{1}{m+1} + \delta_1^2 + \delta_2 \right)}_{\text{III}} \right) \\
& = \frac{1}{\log \frac{m}{q(A+1)k}} (\text{I} + |\text{II}| + |\text{III}|).
\end{aligned}$$

Since $\frac{m}{q(A+1)k} \rightarrow \infty$. The proof would be completed once we show I, |II|, and |III| are bounded. To this end, first note that

$$\begin{aligned}
\text{I} &= \sum_{i=1}^{\lfloor Ak \rfloor} \log^2 \frac{q(k+i)(m+1)}{im} \\
&\leq \sum_{i=1}^{\lfloor Ak \rfloor} \max \left\{ k \int_{(i-1)/k}^{i/k} \log^2 \frac{q(m+1)(1+x)}{mx} dx, k \int_{i/k}^{(i+1)/k} \log^2 \frac{q(m+1)(1+x)}{mx} dx \right\} \\
&\leq 2k \int_0^{A+1} \log^2 \frac{q(m+1)(1+x)}{mx} dx \asymp k = o\left(2k \log \frac{m}{k}\right).
\end{aligned}$$

The second term II obeys

$$\begin{aligned}
|\text{II}| &\leq \sum_{i=1}^{\lfloor Ak \rfloor} 2 \left| \delta_1 \log \frac{q(k+i)(m+1)}{im} \right| \lesssim \sum_{i=1}^{\lfloor Ak \rfloor} \frac{1}{i} \left| \log \frac{q(k+i)(m+1)}{im} \right| \\
&\leq \sum_{i=1}^{\lfloor Ak \rfloor} \max \left\{ k \int_{(i-1)/k}^{i/k} \left| \log \frac{q(m+1)(1+x)}{mx} \right| dx, k \int_{i/k}^{(i+1)/k} \left| \log \frac{q(m+1)(1+x)}{mx} \right| dx \right\} \\
&\leq 2k \int_0^{A+1} \left| \log \frac{q(m+1)(1+x)}{mx} \right| dx \asymp k = o\left(2k \log \frac{m}{k}\right),
\end{aligned}$$

where we use the fact that $\int_0^{A+1} \left| \log \frac{q(m+1)(1+x)}{mx} \right| dx$ is bounded by some constant. The last term is simply bounded as

$$|\text{III}| \leq \sum_{i=1}^{\lfloor Ak \rfloor} \left| \frac{1}{i} - \frac{1}{m+1} + \delta_1^2 + \delta_2 \right| \lesssim \sum_{i=1}^{\lfloor Ak \rfloor} \frac{1}{i} \lesssim \log(Ak) = o\left(2k \log \frac{m}{k}\right).$$

Combining these established bounds on I, II, and III finishes proof. \blacksquare

Lemma E.2. *Under the preceding conditions, let A be any constant satisfying $q(1+A)/A < 1$ and α be sufficiently small such that $l/\lambda_{k+\lfloor \alpha m \rfloor} < 1/2$. Then,*

$$\frac{1}{2k \log(m/k)} \sum_{i=\lfloor Ak \rfloor}^{\lfloor \alpha m \rfloor} \mathbb{E}(\zeta_{(i)} - \lambda_{k+i})_+^2 \rightarrow 0.$$

Proof of Lemma E.2. Note that $\lambda_{k+\lfloor \alpha m \rfloor} \sim \sqrt{2 \log \frac{m}{q(k+\lfloor \alpha m \rfloor)}} \sim \sqrt{2 \log \frac{1}{q\alpha}}$. So it is clear that such α exists. Pick any fixed i between $\lceil Ak \rceil$ and $\lfloor \alpha m \rfloor$. As in the proof of Lemma A.4 in [26], denote by $\alpha_u = \mathbb{P}(\chi_l > \lambda_{k+i} + u)$. Note that

$$\begin{aligned}
\alpha_u &= \mathbb{P}(\chi_l > \lambda_{k+i} + u) = \int_{(\lambda_{k+i} + u)^2}^{\infty} \frac{1}{e^{l/2} \Gamma(l/2)} x^{l/2-1} e^{-x/2} dx \\
&= \int_{\lambda_{k+i}^2}^{\infty} \frac{1}{e^{l/2} \Gamma(l/2)} \left(\frac{(\lambda_{k+i} + u)^2}{\lambda_{k+i}^2} y \right)^{l/2-1} \exp \left(-\frac{(\lambda_{k+i} + u)^2}{2\lambda_{k+i}^2} y \right) d \frac{(\lambda_{k+i} + u)^2}{\lambda_{k+i}^2} y \\
&= \left(1 + \frac{u}{\lambda_{k+i}} \right)^l \int_{\lambda_{k+i}^2}^{\infty} \frac{1}{e^{l/2} \Gamma(l/2)} y^{l/2-1} \exp \left(-\frac{(\lambda_{k+i} + u)^2}{2\lambda_{k+i}^2} y \right) dy \\
&\leq \left(1 + \frac{u}{\lambda_{k+i}} \right)^l e^{-\lambda_{k+i} u} \int_{\lambda_{k+i}^2}^{\infty} \frac{1}{e^{l/2} \Gamma(l/2)} y^{l/2-1} e^{-y/2} dy \\
&= \left(1 + \frac{u}{\lambda_{k+i}} \right)^l e^{-\lambda_{k+i} u} \alpha_0 \\
&\leq \exp \left(\frac{l}{\lambda_{k+i}} u - \lambda_{k+i} u \right) \alpha_0.
\end{aligned}$$

With the proviso that $l/\lambda_{k+\lfloor \alpha m \rfloor} < 1/2 < \lambda_{k+\lfloor \alpha m \rfloor}/2$, it follows that

$$\alpha_u \leq e^{-\lambda_{k+i} u/2} \alpha_0.$$

The remaining proof follows from exactly the same reasoning as that of Lemma A.4 in [26]. ■

Lemma E.3. *Under the preceding conditions, for any constant $\alpha > 0$ we have*

$$\frac{1}{2k \log(m/k)} \sum_{i=\lfloor \alpha m \rfloor}^{m-k} \mathbb{E}(\zeta_{(i)} - \lambda_{k+i})_+^2 \rightarrow 0.$$

Proof of Lemma E.3. Recognizing that the value of the summation increases as α decreases, we only prove the lemma for sufficiently small α . In the case of $U_{(i)} \geq \alpha/3$, we get

$$\begin{aligned}
(\zeta_{(i)} - \lambda_{k+i})_+ &= (F_{\chi_l}^{-1}(1 - U_{(i)}) - F_{\chi_l}^{-1}(1 - q(k+i)/m))_+ \\
&\asymp (1 - U_{(i)} - (1 - q(k+i)/m))_+ \\
&= (q(k+i)/m - U_{(i)})_+,
\end{aligned}$$

since both $U_{(i)}$ and $q(k+i)/m$ are bounded below away from zero. Otherwise, we use the trivial inequality $(\zeta_{(i)} - \lambda_{k+i})_+ \leq \zeta_{(i)}$. In either case, we get

$$\begin{aligned}
(\zeta_{(i)} - \lambda_{k+i})_+^2 &\lesssim \zeta_{(i)}^2 \mathbf{1}_{U_{(i)} < \frac{\alpha}{3}} + \left(\frac{q(k+i)}{m} - U_{(i)} \right)_+^2 \\
&= \left(F_{\chi_l}^{-1}(1 - U_{(i)}) \right)^2 \mathbf{1}_{U_{(i)} < \frac{\alpha}{3}} + \left(\frac{q(k+i)}{m} - U_{(i)} \right)_+^2 \\
&\asymp 2 \log \left(\frac{1}{U_{(i)}} \right) \mathbf{1}_{U_{(i)} < \frac{\alpha}{3}} + \left(\frac{q(k+i)}{m} - U_{(i)} \right)_+^2 \\
&\lesssim \log \left(\frac{1}{U_{(i)}} \right) \mathbf{1}_{U_{(i)} < \frac{\alpha}{3}} + \mathbf{1}_{U_{(i)} \leq \frac{q(k+i)}{m}}.
\end{aligned}$$

Hence,

$$\sum_{i=\lceil \alpha m \rceil}^{m-k} \mathbb{E}(\zeta_{(i)} - \lambda_{k+i})_+^2 \lesssim \sum_{i=\lceil \alpha m \rceil}^{m-k} \mathbb{E} \left(\log \left(\frac{1}{U_{(i)}} \right); U_{(i)} < \frac{\alpha}{3} \right) + \sum_{i=\lceil \alpha m \rceil}^{m-k} \mathbb{P} \left(U_{(i)} \leq \frac{q(k+i)}{m} \right)$$

In the remaining proof we aim to show

$$\sum_{i=\lceil \alpha m \rceil}^{m-k} \mathbb{E} \left(\log \left(\frac{1}{U_{(i)}} \right); U_{(i)} < \frac{\alpha}{3} \right) \rightarrow 0 \quad (\text{E.6})$$

and

$$\sum_{i=\lceil \alpha m \rceil}^{m-k} \mathbb{P} \left(U_{(i)} \leq \frac{q(k+i)}{m} \right) \rightarrow 0. \quad (\text{E.7})$$

This is more than we need since $2k \log(m/k) \rightarrow \infty$.

Each summand of (E.6) is bounded above by

$$\begin{aligned} \mathbb{E} \left(\log \left(\frac{1}{U_{(\lceil \alpha m \rceil)}} \right); U_{(\lceil \alpha m \rceil)} < \frac{\alpha}{3} \right) &= \int_0^{\frac{\alpha}{3}} \frac{x^{\lceil \alpha m \rceil - 1} (1-x)^{m - \lceil \alpha m \rceil} \log \frac{1}{x}}{\text{B}(\lceil \alpha m \rceil, m + 1 - \lceil \alpha m \rceil)} dx \\ &\leq \int_0^{\frac{\alpha}{3}} \frac{x^{\lceil \alpha m \rceil - 1} \log \frac{1}{x}}{\text{B}(\lceil \alpha m \rceil, m + 1 - \lceil \alpha m \rceil)} dx \\ &= \frac{1}{\lceil \alpha m \rceil^2 \text{B}(\lceil \alpha m \rceil, m + 1 - \lceil \alpha m \rceil)} \int_0^{(\frac{\alpha}{3})^{\lceil \alpha m \rceil}} \log \frac{1}{y} dy \\ &\sim \frac{(\alpha/3)^{\lceil \alpha m \rceil} \log \frac{3}{\alpha}}{\lceil \alpha m \rceil \text{B}(\lceil \alpha m \rceil, m + 1 - \lceil \alpha m \rceil)}. \end{aligned}$$

The last line obeys

$$\begin{aligned} \log \left[\frac{(\alpha/3)^{\lceil \alpha m \rceil}}{\text{B}(\lceil \alpha m \rceil, m + 1 - \lceil \alpha m \rceil)} \right] &\sim -\alpha m \log \frac{3}{\alpha} + \alpha m \log \frac{1}{\alpha} + (1 - \alpha)m \log \frac{1}{1 - \alpha} \\ &= -\alpha m \log 3 + (1 - \alpha)m \log \frac{1}{1 - \alpha}. \end{aligned}$$

For small α , we get $-\alpha \log 3 + (1 - \alpha) \log \frac{1}{1 - \alpha} = -\alpha \log 3 + (1 + o(1))(1 - \alpha)\alpha = -(\log 3 - 1 + o(1))\alpha$. (Note that $\log 3 - 1 = 0.0986 \dots > 0$.) This immediately yields

$$\mathbb{E} \left(\log \left(\frac{1}{U_{(\lceil \alpha m \rceil)}} \right); U_{(\lceil \alpha m \rceil)} < \frac{\alpha}{3} \right) \sim e^{-(\log 3 - 1 + o(1))\alpha m},$$

which implies (E.6) since $m e^{-(\log 3 - 1 + o(1))\alpha m} \rightarrow 0$.

Next, we turn to show (E.7). Note that $\mathbb{P} \left(U_{(i)} \leq \frac{q(k+i)}{m} \right)$ actually is the tail probability of the binomial distribution with m trials and success probability $\frac{q(k+i)}{m}$. Hence, by the Chernoff bound, this probability is bounded as

$$\mathbb{P} \left(U_{(i)} \leq \frac{q(k+i)}{m} \right) \leq \exp(-m \text{KL}(i/m || q(k+i)/m)),$$

where $\text{KL}(a||b) := a \log \frac{a}{b} + (1 - a) \log \frac{1-a}{1-b}$ is the Kullback-Leibler divergence. Thanks to $i \geq \lceil \alpha m \rceil \gg k$, simple analysis reveals that

$$\text{KL}(i/m || q(k+i)/m) \geq (1 + o(1))i \left(\log \frac{1}{q} - 1 + q \right) / m.$$

Combining the last two displays gives

$$\mathbb{P}\left(U_{(i)} \leq \frac{q(k+i)}{m}\right) \leq e^{-(1+o(1))(\log \frac{1}{q}-1+q)i}.$$

Plugging the above inequality into (E.7) yields

$$\sum_{i=\lceil \alpha m \rceil}^{m-k} \mathbb{P}\left(U_{(i)} \leq \frac{q(k+i)}{m}\right) \leq \sum_{i=\lceil \alpha m \rceil}^{m-k} e^{-(1+o(1))(\log \frac{1}{q}-1+q)i} \rightarrow 0,$$

where the last step follows from $\log \frac{1}{q} - 1 + q > 0$ and $\lceil \alpha m \rceil \rightarrow \infty$. \blacksquare

F Strength of signals

Consider the case when all submatrices X_{I_i} have the same rank, $l > 0$, $w > 0$ is used as the universal weight and X is orthogonal at groups level. From the interpretation of gSLOPE estimate coming from (2.8), we see that the identification of the relevant groups could be summarized as follows: λ decides on the number, R , of groups labeled as relevant, which correspond to indices of the R largest values among $w^{-1}\|\tilde{y}_{I_1}\|_2, \dots, w^{-1}\|\tilde{y}_{I_m}\|_2$. The random variables $w^{-1}\|\tilde{y}_{I_i}\|_2$ have a (possibly) non-central χ distributions with l degrees of freedom and noncentrality parameters given by the entries of $\|\tilde{\beta}\|_{\mathbb{I}}$. Now, the nonzero $\|\tilde{\beta}_{I_i}\|_2$ could be perceived as a strong signal, if with the high probability the random variable having the noncentral χ distribution with the noncentrality parameter $\|\tilde{\beta}_{I_i}\|_2$ is large comparing to the background composed of the independent random variables with the χ_l distributions (then signal is likely to be identified by gSLOPE; otherwise, the signal could be easily covered by random disturbances and its identification has more in common with good luck than with the usage of particular method). The important quantity, which could be treated as a breaking point, is the expected value of the maximum of the background noise. Group effects being close to this value, could be perceived as medium under the orthogonal case and weak under the occurrence of correlations between groups. The above reasoning applied to the considered case, yields the issue of approximation of the expected value of the maximum of m independent χ_l -distributed variables. Suppose that $\Psi_i \sim \chi_l$ for $i = \{1, \dots, m\}$. From Jensen's inequality we have

$$\mathbb{E}\left(\max_{i=1, \dots, m} \{\Psi_i\}\right) = \mathbb{E}\left(\sqrt{\max_{i=1, \dots, m} \{\Psi_i^2\}}\right) \leq \sqrt{\mathbb{E}\left(\max_{i=1, \dots, m} \{\Psi_i^2\}\right)},$$

hence we will replace the last problem by the problem of finding the reasonable upper bound on the expected value of the maximum of m independent, χ_l^2 -distributed variables.

Theorem F.1. *Let Ψ_1, \dots, Ψ_m be independent variables, $\Psi_i \sim \chi_l^2$ for all i . Then*

$$\mathbb{E}\left(\max_{i=1, \dots, m} \{\Psi_i\}\right) \leq \frac{4 \ln(m)}{1 - m^{-\frac{2}{l}}}. \quad (\text{F.1})$$

Proof. Denote $M_m := \max_{i=1, \dots, m} \{\Psi_i\}$. From the Jensen's inequality applied to e^{tM_m} we have

$$e^{t\mathbb{E}[M_m]} \leq \mathbb{E}[e^{tM_m}] = \mathbb{E}\left[\max_{i=1, \dots, m} e^{t\Psi_i}\right] \leq \sum_{i=1}^m \mathbb{E}[e^{t\Psi_i}]. \quad (\text{F.2})$$

We will consider only $t \in [0, \frac{1}{2})$. Since the moment generating function for χ_l^2 distribution is given by $MGF := (1-2t)^{-\frac{l}{2}}$, for each i it holds $\mathbb{E}[e^{t\Psi_i}] = (1-2t)^{-\frac{l}{2}}$ and we get $e^{t\mathbb{E}[M_m]} \leq m(1-2t)^{-\frac{l}{2}}$. Applying the natural logarithm to both sides yields

$$\mathbb{E}[M_m] \leq \frac{\ln(m) + \ln\left((1-2t)^{-\frac{l}{2}}\right)}{t}, \quad t \in [0, 1/2). \quad (\text{F.3})$$

Define $t_{m,l} := \frac{1-m^{-\frac{2}{l}}}{2}$. Then for all positive, natural numbers l and m we have $t_{m,l} \in [0, \frac{1}{2}]$. Plugging $t_{m,l}$ to the right side of (F.3) gives inequality (F.1) and finishes the proof. \blacksquare

The above theorem gives us the motivation to use the quantity $\sqrt{4 \ln(m)/(1 - m^{-2/l})}$ as the upper bound on the expected value of maximum over m independent χ_l -distributed variables. In all simulations, which we have performed to investigate the performance of gSLOPE, we have generated the effects for truly relevant groups basing on these upper bounds. In particular, in experiments where l_i 's as well as weights were identical, we aimed at $\mathbb{E}(\|\tilde{y}_{\mathbb{I}_i}\|_2) = \sqrt{4 \ln(m)/(1 - m^{-2/l})}$, for the truly relevant group i . Since $\mathbb{E}(\|\tilde{y}_{\mathbb{I}_i}\|_2) \approx \sqrt{\|\tilde{\beta}_{\mathbb{I}_i}\|_2^2 + l}$, this yields the setting

$$\|\tilde{\beta}_{\mathbb{I}_i}\|_2 = B(m, l), \quad \text{for} \quad B(m, l) := \sqrt{4 \ln(m)/(1 - m^{-2/l}) - l} \quad (\text{F.4})$$

for groups chosen to be truly relevant.

G The sequence of tuning parameters when variables in different groups are independent

To model the situation when variables in different groups are stochastically independent we will assume that n by p design matrix is a realization of the random matrix with independent entries drawn from the normal distribution, $\mathcal{N}(0, \frac{1}{n})$, so as the expected value of $X_i^\top X_j$ is equal to 1 for $i = j$, and equal to 0 otherwise. The main objective is to derive the lambda sequence, which could be applied to achieve gFDR control under assumption that the $\llbracket \beta \rrbracket_{X,I}$ is sparse. At first we will confine ourselves only to the case $l_1 = \dots = l_m := l$, $w_1 = \dots = w_m := w$ and when the number of elements in each group is relatively small as compared to the number of observations ($l \ll n$). For simplicity in this subsection we will fix $\sigma = 1$. In case when $\sigma \neq 1$, the proposed sequence lambda should be multiplied by σ , as in expression (2.2). In the heuristics presented in this subsection, we will use the notation $A \approx B$, in order to express that with large probability the differences between corresponding entries of matrices A and B are very small.

In situation when entries of X come from $\mathcal{N}(0, \frac{1}{n})$ distribution and sizes of groups are relatively small, a very good approximation of β^{gs} could be obtained by $\hat{\beta}$, defined as

$$\hat{\beta} := \arg \min_b \left\{ \frac{1}{2} \|y - Xb\|_2^2 + \sigma J_\lambda(W \llbracket b \rrbracket_I) \right\}. \quad (\text{G.1})$$

Assume for simplicity that $\|\beta_{I_1}\|_2 > \dots > \|\beta_{I_s}\|_2 > 0$, $\|\beta_{I_j}\|_2 = 0$ for $j > s$, $\hat{\beta}$ satisfies the same conditions for some λ and the true model is sparse. Divide I into two families of sets $I^s := \{I_1, \dots, I_s\}$ and $I^c := \{I_{s+1}, \dots, I_m\}$. To derive optimality condition for $\hat{\beta}$ we will prove the following

Theorem G.1. *Let $b \in \mathbb{R}^p$ be such that $\|b_{I_1}\|_2 > \dots > \|b_{I_s}\|_2 > 0$, $\|b_{I_j}\|_2 = 0$ for $j > s$ and denote $\lambda^c := (\lambda_{s+1}, \dots, \lambda_m)^\top$. If $g \in \partial J_\lambda(w \llbracket b \rrbracket_I)$, then it holds:*

$$\begin{cases} g_{I_i} = w \lambda_i \frac{b_{I_i}}{\|b_{I_i}\|_2}, \quad i = 1, \dots, s \\ \llbracket g \rrbracket_{I^c} \in C_{w\lambda^c} \end{cases}, \quad (\text{G.2})$$

where the set C_λ (here with $w\lambda^c$ instead of λ) is defined in appendix (B.4).

Proof. For $b \in \mathbb{R}^p$ define $J_{\lambda,I}(b) := J_\lambda(\llbracket b \rrbracket_I)$ and put $H := \{h \in \mathbb{R}^p : \|(b+h)_{I_1}\|_2 > \dots > \|(b+h)_{I_s}\|_2, \|(b+h)_{I_s}\|_2 > \|h_{I_j}\|_2, j > s\}$. If $g \in \partial J_{\lambda,I}(b)$, then for all $h \in H$ from definition of subgradient it holds

$$\sum_{i=1}^s \lambda_i \|(b+h)_{I_i}\|_2 + \sum_{i=s+1}^m \lambda_i (\llbracket b+h \rrbracket_I)_{(i)} \geq \sum_{i=1}^s \lambda_i \|b_{I_i}\|_2 + \sum_{i=1}^s g_{I_i}^\top h_{I_i} + (g^c)^\top h^c, \quad (\text{G.3})$$

for $g^c := (g_{I_{s+1}}^\top, \dots, g_{I_m}^\top)^\top$ and $h^c := (h_{I_{s+1}}^\top, \dots, h_{I_m}^\top)^\top$. Define $\tilde{I} := \{\tilde{I}_1, \dots, \tilde{I}_{m-s}\}$, with set $\tilde{I}_i := \{(i-1) \cdot l + 1, \dots, i \cdot l\}$. Then $\llbracket g^c \rrbracket_{\tilde{I}} = \llbracket g \rrbracket_{I^c}$. Consider first case, when h belongs to the set $H^c := \{h \in H : h_{I_i} \equiv 0, i \leq s\}$. This yields

$$\sum_{i=1}^{m-s} \lambda_{s+i} (\llbracket h^c \rrbracket_{\tilde{I}})_{(i)} \geq (g^c)^\top h^c. \quad (\text{G.4})$$

Since $\{h^c : h \in H^c\}$ is open in $\mathbb{R}^{l(m-s)}$ and contains zero, from Proposition G.5 we have that $g^c \in \partial J_{\lambda^c, \tilde{I}}(0)$ and the inequality (G.4) is true for any $h^c \in \mathbb{R}^{l(m-s)}$ yielding

$$0 \geq \sup_{h^c} \left\{ (g^c)^\top h^c - J_{\lambda^c, \tilde{I}}(h^c) \right\} = J_{\lambda^c, \tilde{I}}^*(g^c) = \begin{cases} 0, & \llbracket g^c \rrbracket_{\tilde{I}} \in C_{\lambda^c} \\ \infty, & \text{otherwise} \end{cases}, \quad (\text{G.5})$$

see Proposition B.3. This result immediately gives condition $\llbracket g^c \rrbracket_{\tilde{I}} \in C_{\lambda^c}$, which is equivalent with $\llbracket g \rrbracket_{I^c} \in C_{\lambda^c}$. To find conditions for g_{I_i} with $i \leq s$, define sets $H_i := \{h \in H : h_{I_j} \equiv 0, j \neq i\}$. For $h \in H_i$, (G.3) reduces to $\lambda_i \|b_{I_i} + h_{I_i}\|_2 \geq \lambda_i \|b_{I_i}\|_2 + g_{I_i}^\top h_{I_i}$. Since the set $\{h_{I_i} : h \in H_i\}$ is open in \mathbb{R}^l and contains zero, from Proposition G.5 we have $g_{I_i} \in \partial f_i(b_{I_i})$ for $f_i : \mathbb{R}^l \rightarrow \mathbb{R}$, $f_i(x) := \lambda_i \|x\|_2$. Since f_i is convex and differentiable in b_{I_i} , it holds $g_{I_i} = \lambda_i \frac{b_{I_i}}{\|b_{I_i}\|_2}$, which finishes the proof. \blacksquare

The above theorem allows to write the optimality condition for $\hat{\beta}$ in form

$$\begin{cases} X_{I_i}^\top (y - X\hat{\beta}) = w\lambda_i \frac{\hat{\beta}_{I_i}}{\|\hat{\beta}_{I_i}\|_2}, & i = 1, \dots, s \\ \llbracket X^\top (y - X\hat{\beta}) \rrbracket_{I^c} \in C_{w\lambda^c} \end{cases}. \quad (\text{G.6})$$

Since $X_{I_i}^\top X_{I_i} \approx \mathbf{I}_i$, for $i \leq s$ we get $X_{I_i}^\top (y - X_{\setminus I_i} \hat{\beta}_{\setminus I_i}) \approx \hat{\beta}_{I_i} \left(1 + \frac{w\lambda_i}{\|\hat{\beta}_{I_i}\|_2}\right)$, where $X_{\setminus I_i}$ is matrix X without columns from I_i and $\hat{\beta}_{\setminus I_i}$ denotes vector $\hat{\beta}$ with removed coefficients indexed by I_i . This means that, for $i = 1, \dots, s$, vector $v_{I_i} := X_{I_i}^\top (y - X_{\setminus I_i} \hat{\beta}_{\setminus I_i})$ is approximately collinear with $\hat{\beta}_{I_i}$. Since $1 + \frac{w\lambda_i}{\|\hat{\beta}_{I_i}\|_2} > 0$, we have $\frac{v_{I_i}}{\|v_{I_i}\|_2} \approx \frac{\hat{\beta}_{I_i}}{\|\hat{\beta}_{I_i}\|_2}$. This yields $\hat{\beta}_{I_i} \approx \left(1 - \frac{w\lambda_i}{\|v_{I_i}\|_2}\right) v_{I_i}$ and consequently $\|\hat{\beta}_{I_i}\|_2 \approx \left| \|v_{I_i}\|_2 - w\lambda_i \right|$. Therefore (G.6) can be written as

$$\begin{cases} \left| \|v_{I_i}\|_2 - w\lambda_i \right| \approx \|\hat{\beta}_{I_i}\|_2, & i = 1, \dots, s \\ \llbracket v \rrbracket_{I^c} \in C_{w\lambda^c} \end{cases}, \quad (\text{G.7})$$

for $v := (v_{I_1}^\top, \dots, v_{I_m}^\top)^\top$.

The task now is to select λ_i 's such that condition $\llbracket v \rrbracket_{I^c} \in C_{w\lambda^c}$ regulates the rate of false discoveries. Denote $I_S := \bigcup_{i=1}^s I_i$. Putting $y = X_{I_S} \beta_{I_S} + z$, we obtain

$$v_{I_i} = X_{I_i}^\top X_{I_S} (\beta_{I_S} - \hat{\beta}_{I_S}) + X_{I_i}^\top z, \quad (\text{G.8})$$

for $i > s$ (irrelevant groups). Under orthogonal design this expression reduces only to the term $X_{I_i}^\top z$, and in such situation $\|v_{I_i}\|_2$ has χ distribution with l degrees of freedom which was used in subsection G to define the sequence λ . In the considered near-orthogonal situation, the term $X_{I_i}^\top X_{I_S} (\beta_{I_S} - \hat{\beta}_{I_S})$ should be also taken into account. The following two assumptions will be important to derive the appropriate approximation of v_{I_i} distribution:

- the distribution of v_{I_i} could be well approximated by multivariate normal distribution,
- for relatively strong effects it occurs $\frac{\hat{\beta}_{I_i}}{\|\hat{\beta}_{I_i}\|_2} \approx \frac{\beta_{I_i}}{\|\beta_{I_i}\|_2}$ for $i = 1, \dots, s$.

The first assumption is justified when one works with large data scenario, based on the Central Limit Theorem. In discussion concerning the second assumption it is important to clarify the effect of penalty imposed on entire groups. The magnitudes of coefficients in $\hat{\beta}_{I_i}$, for truly relevant group i , are generally significantly smaller than in β_{I_i} . This, a so-called shrinking effect, is typical for penalized methods. It turns out, however, that under assumed conditions estimates of coefficients of nonzero β_{I_i} are pulled to zero proportionally and after normalizing, $\hat{\beta}_{I_i}$ and β_{I_i} are comparable.

From the upper equation in (G.6), we have that $X_{I_s}^\top(X_{I_s}\beta_{I_s} - X_{I_s}\hat{\beta}_{I_s}) + X_{I_s}^\top z \approx wH_{\lambda,\beta}$, for

$$H_{\lambda,\beta} := \left(\lambda_1 \frac{\beta_{I_1}^\top}{\|\beta_{I_1}\|_2}, \dots, \lambda_s \frac{\beta_{I_s}^\top}{\|\beta_{I_s}\|_2} \right)^\top, \quad (\text{G.9})$$

which gives $X_{I_i}^\top X_{I_s}(\beta_{I_s} - \hat{\beta}_{I_s}) \approx X_{I_i}^\top X_{I_s}(X_{I_s}^\top X_{I_s})^{-1}(wH_{\lambda,\beta} - X_{I_s}^\top z)$. Combining the last expression with (G.8) yields

$$v_{I_i} \approx X_{I_i}^\top X_{I_s}(X_{I_s}^\top X_{I_s})^{-1}(wH_{\lambda,\beta} - X_{I_s}^\top z) + X_{I_i}^\top z. \quad (\text{G.10})$$

To determine the parameters of multivariate normal distribution, which best describes the distribution of v_{I_i} , we will derive the exact values of the mean and the covariance matrix of the distribution of the right-hand side expression in (G.10) for $i > s$. Since $I_i \cap I_s = \emptyset$ and entries of X matrix are randomized independently with $\mathcal{N}(0, \frac{1}{n})$ distribution, the expected value of the random vector in (G.10) is 0 and its covariance matrix is provided by the following Lemma.

Lemma G.2. *The covariance matrix of $\hat{v}_{I_i} := X_{I_i}^\top X_{I_s}(X_{I_s}^\top X_{I_s})^{-1}(wH_{\lambda,\beta} - X_{I_s}^\top z) + X_{I_i}^\top z$, for $i > s$, is given by the formula*

$$\text{Cov}(\hat{v}_{I_i}) = \left(\frac{n-ls}{n} + w^2 \frac{\|\lambda^S\|_2^2}{n-ls-1} \right) \mathbf{I}_l,$$

where $\lambda^S := (\lambda_1, \dots, \lambda_s)^\top$.

Before proving Lemma G.2, we will introduce two auxiliary results, proofs of which can be found at the end of this section.

Lemma G.3. *Suppose that entries of a random matrix $X \in M(n, r)$, with $r \leq n$, are independently and identically distributed and have a normal distribution with zero mean. Then, there exists the expected value of a random matrix $A_X = X(X^\top X)^{-1}X^\top$ and $\mathbb{E}(A_X) = \frac{r}{n}\mathbf{I}_n$.*

Lemma G.4. *Suppose that $X \in M(n, r)$, with $r+1 < n$, and entries of X are independent and identically distributed, $X_{ij} \sim \mathcal{N}(0, 1/n)$ for all i and j . Then, there exists expected value of random matrix, $M_{X,\lambda} := B_X H_{\lambda,\beta} H_{\lambda,\beta}^\top B_X^\top$, for $B_X = X(X^\top X)^{-1}$ and $H_{\lambda,\beta}$ defined in (G.9). Moreover, it holds $\mathbb{E}(M_{X,\lambda}) = \frac{\|\lambda^S\|_2^2}{n-r-1}\mathbf{I}_n$.*

Proof of Lemma G.2. We have $\hat{v}_{I_i} = \xi_{X,z} + \zeta_X$, for $\xi_{X,z} := X_{I_i}^\top (\mathbf{I}_n - A_X)z$, $\zeta_X := wX_{I_i}^\top B_X H_{\lambda,\beta}$, $A_X := X_{I_s}(X_{I_s}^\top X_{I_s})^{-1}X_{I_s}^\top$, $B_X := X_{I_s}(X_{I_s}^\top X_{I_s})^{-1}$. Since $\mathbb{E}(\xi_{X,z}\zeta_X^\top) = 0$ and mean of \hat{v}_{I_i} is equal to 0, it holds $\text{Cov}(\hat{v}_{I_i}) = \text{Cov}(\xi_{X,z}) + \text{Cov}(\zeta_X)$. Now thanks to Lemma G.3 and Lemma G.4

$$\begin{aligned} \text{Cov}(\xi_{X,z}) &= \mathbb{E} \left[X_{I_i}^\top (\mathbf{I}_n - A_X) z z^\top (\mathbf{I}_n - A_X)^\top X_{I_i} \right] = \\ &= \mathbb{E} \left[X_{I_i}^\top (\mathbf{I}_n - A_X) (\mathbf{I}_n - A_X)^\top X_{I_i} \right] = \mathbb{E} \left[X_{I_i}^\top (\mathbf{I}_n - A_X) X_{I_i} \right] = \\ &= \frac{1}{n} (n-ls) \cdot \mathbb{E} \left[X_{I_i}^\top X_{I_i} \right] = \frac{1}{n} (n-ls) \cdot \mathbf{I}_l, \end{aligned} \quad (\text{G.11})$$

$$\begin{aligned} \text{Cov}(\zeta_X) &= w^2 \mathbb{E} \left[X_{I_i}^\top B_X H_{\lambda,\beta} H_{\lambda,\beta}^\top B_X^\top X_{I_i} \right] = w^2 \frac{\|\lambda^S\|_2^2}{n-sl-1} \mathbb{E} \left[X_{I_i}^\top X_{I_i} \right] = \\ &= w^2 \frac{\|\lambda^S\|_2^2}{n-sl-1} \mathbf{I}_l, \end{aligned} \quad (\text{G.12})$$

which finishes the proof. ■

We have shown that for $i > s$ the distribution of $\|v_{I_i}\|_2$ could be approximated by scaled χ distribution with l degrees of freedom and scale parameter $\mathcal{S} = \sqrt{\frac{n-ls}{n} + \frac{w^2\|\lambda_S\|_2^2}{n-sl-1}}$. Now, analogously to the orthogonal situation, lambdas could be defined as $\lambda_i := \frac{1}{w_i} F_{\mathcal{S}\chi_l}^{-1} \left(1 - \frac{q^i}{m}\right) = \frac{\mathcal{S}}{w_i} F_{\chi_l}^{-1} \left(1 - \frac{q^i}{m}\right)$. Since s is unknown, we will apply the strategy used in [10]: define λ_1 as in orthogonal case and for $j \geq 2$ define λ_j basing on already generated sequence, according to following procedure.

Procedure 6 Selecting lambdas in near-orthogonal situation: equal groups sizes

input: $q \in (0, 1)$, $w > 0$, $p, n, m, l \in \mathbb{N}$

$\lambda_1 := \frac{1}{w} F_{\chi_l}^{-1} \left(1 - \frac{q}{m}\right)$;

For $i \in \{2, \dots, m\}$:

$\lambda^S := (\lambda_1, \dots, \lambda_{i-1})^\top$;

$\mathcal{S} := \sqrt{\frac{n-l(i-1)}{n} + \frac{w^2\|\lambda^S\|_2^2}{n-l(i-1)-1}}$;

$\lambda_i^* := \frac{\mathcal{S}}{w} F_{\chi_l}^{-1} \left(1 - \frac{q^i}{m}\right)$;

if $\lambda_i^* \leq \lambda_{i-1}$, then put $\lambda_i := \lambda_i^*$. Otherwise, stop the procedure and put $\lambda_j := \lambda_{i-1}$ for $j \geq i$;

end for

Consider now the Gaussian design with arbitrary group sizes and sequence of positive weights w_1, \dots, w_m . One possible approach is to construct consecutive λ_i as the largest scaled quantiles among all distributions, i.e. as $\max_{j=1, \dots, m} \left\{ \frac{\mathcal{S}_j}{w_j} F_{\chi_{l_j}}^{-1} \left(1 - \frac{q^i}{m}\right) \right\}$ for corrections \mathcal{S}_j 's adjusted to different l_i values (the conservative strategy). In this article, however, we will stick to the more liberal strategy based on λ^{mean} , which leads to the modified sequence of tuning parameters presented in Procedure 1.

Proposition G.5. *For any open set H containing zero the subgradient of convex function f at b could be equivalently defined as a vector g satisfying $f(b+h) \geq f(b) + g^\top h$, for all $h \in H$.*

Proof. Suppose that f is convex function and for some $b, g \in \mathbb{R}^p$ it occurs $f(b+h) \geq f(b) + g^\top h$ for $h \in H$, where H is open set containing zero. Let $h_0 \in \mathbb{R}^p$ be arbitrary vector. Function $F: \mathbb{R} \rightarrow \mathbb{R}$, defined as $F(t) := f(b+th_0) - tg^\top h_0$, is convex. There exists $t_0 \in (0, 1)$ such that $t_0 h_0 \in H$, what gives

$$f(b) \leq F(t_0) = F((1-t_0) \cdot 0 + t_0 \cdot 1) \leq (1-t_0)f(b) + t_0 F(1) \quad (\text{G.13})$$

and $f(b+h_0) \geq f(b) + g^\top h_0$ as a result. \blacksquare

G.1 The proof of Lemma G.3

The claim is obvious for $n = 1$ and we will assume that $n > 1$. First, we will list some basic properties of A_X . It could be easily noticed that: A_X is symmetric matrix, A_X is idempotent matrix (meaning that $A_X A_X = A_X$) and that $\text{trace}(A_X) = \text{trace}(X^\top X (X^\top X)^{-1}) = r$. We will now show that for each $i \in \{1, \dots, n\}$, $j \in \{1, \dots, r\}$ the support of a $A_X(i, j)$ distribution is bounded, which will give us the existence of the expected value. Let $\|A\|_F$ be the Frobenius norm. Then

$$|(A_X)_{i,j}| \leq \|A\|_F = \sqrt{\text{trace}(A_X^\top A_X)} = \sqrt{\text{trace}(A_X)} = \sqrt{r}. \quad (\text{G.14})$$

We will use notation $E_X := \mathbb{E}(A_X)$. Since entries of matrix X are randomized independently with the same distribution, E_X is invariant under permutation applied to rows, i.e. $E_X = E_{PX}$ for any permutation matrix P . This gives $E_X = P E_X P^\top$, which means that applying the same permutation to rows and columns has no impact on expected value. We will show that

$$(E_X)_{i,j} = (E_X)_{1,n}, \text{ for } i < j. \quad (\text{G.15})$$

Consider first the case when $i = 1$ and $1 < j < n$. Denoting by $P_{j \leftrightarrow n}$ matrix corresponding to transposition which replaces elements j and n , we have $(E_X)_{1,j} = (P_{j \leftrightarrow n} E_X P_{j \leftrightarrow n}^\top)_{1,j} = (E_X)_{1,n}$.

When $j = n$ and $1 < i < n$, the same reasoning works with $P_{1 \leftrightarrow i}$. Suppose now, that $1 < i < n$ and $1 < j < n$. We get $(E_X)_{i,j} = (E_X)_{1,n}$ analogously by using arbitrary permutation matrix P which replaces element j with n and element i with 1. Since A_X is symmetric, (G.15) is true also for $i > j$. On the other hand, for all $i, j \in \{1, \dots, n\}$, we have $(E_X)_{i,i} = (P_{j \leftrightarrow i} E_X P_{j \leftrightarrow i}^\top)_{i,i} = (E_X)_{j,j}$. Consequently, all off-diagonal entries of E_X are equal to some t and all diagonal entries have the same value d . Since

$$nd = \text{trace}(E_X) = \sum_{i=1}^n \mathbb{E}(A_X(i, i)) = \mathbb{E} \left(\sum_{i=1}^n (A_X)_{i,i} \right) = r, \quad (\text{G.16})$$

we have $d = \frac{r}{n}$ and it remains to show that $t = 0$. Define $\Sigma := \begin{bmatrix} -1 & \mathbf{0}^\top \\ \mathbf{0} & \mathbf{I}_{n-1} \end{bmatrix}$. Then ΣX_S differs from X_S only by signs of the first row. Since entries of matrix X_S have zero-symmetric distribution, we have $E_X = E_{\Sigma X}$. Now

$$\begin{bmatrix} d & \mathbf{1}_{n-1}^\top t \\ \mathbf{1}_{n-1} t & \ddots \end{bmatrix} = E_X = \Sigma E_X \Sigma = \begin{bmatrix} d & -\mathbf{1}_{n-1}^\top t \\ -\mathbf{1}_{n-1} t & \ddots \end{bmatrix}, \quad (\text{G.17})$$

which implies $t = 0$ and proves the statement.

G.2 The proof of Lemma G.4

It is easy to see that $M_{X,\lambda}$ is symmetric, positive semi-definite matrix. Denote by $\|M_{X,\lambda}\|_*$ the nuclear (trace) norm of matrix $M_{X,\lambda}$. We have

$$\begin{aligned} \mathbb{E}|(M_{X,\lambda})_{i,j}| &\leq \mathbb{E}(\|M_{X,\lambda}\|_*) = \mathbb{E}(\text{trace}(M_{X,\lambda})) = \mathbb{E}(\text{trace}(H_{\lambda,\beta}^\top B_X^\top B_X H_{\lambda,\beta})) = \\ \mathbb{E}(H_{\lambda,\beta}^\top (X^\top X)^{-1} H_{\lambda,\beta}) &= \frac{n}{(n-r-1)} H_{\lambda,\beta}^\top H_{\lambda,\beta} = \frac{n \|\lambda^S\|_2^2}{n-r-1}, \end{aligned} \quad (\text{G.18})$$

since $X^\top X$ follows an inverse Wishart distribution. This gives the existence of $E_X := \mathbb{E}(M_{X,\lambda})$. Analogously to situation in Lemma G.3, E_X is invariant under permutation or signs changes applied to rows of X , i.e. $E_X = E_{PX}$ for any permutation matrix P , and $E_X = E_{\Sigma X}$ for diagonal matrix Σ with entries on diagonal coming from set $\{-1, 1\}$. Since $E_{PX} = PE_X P^\top$ and $E_{\Sigma X} = \Sigma E_X \Sigma$, as before we have that E_X is diagonal matrix with all diagonal entries having the same value d . The value d could be easily found using (G.18) since we have

$$nd = \text{trace}(E_X) = \frac{n \|\lambda^S\|_2^2}{n-r-1}. \quad (\text{G.19})$$