

1 Spatially relaxed inference
2 on high-dimensional linear models

3 Jérôme-Alexis CHEVALIER
 Inria Paris-Saclay, CEA, Université Paris-Saclay
 and
 Tuan-Binh NGUYEN
 Inria Paris-Saclay, CEA, Université Paris-Saclay, LMO
 and
 Bertrand THIRION
 Inria Paris-Saclay, CEA, Université Paris-Saclay
 and
 Joseph SALMON
 IMAG, Université de Montpellier, CNRS

jerome-alexis.chevalier@inria.fr

4 June 1, 2021

5 **Abstract**

6 We consider the inference problem for high-dimensional linear models, when co-
7 variates have an underlying spatial organization reflected in their correlation. A
8 typical example of such a setting is high-resolution imaging, in which neighboring
9 pixels are usually very similar. Accurate point and confidence intervals estimation
10 is not possible in this context with many more covariates than samples, furthermore
11 with high correlation between covariates. This calls for a reformulation of the sta-
12 tistical inference problem, that takes into account the underlying spatial structure:
13 if covariates are locally correlated, it is acceptable to detect them up to a given
14 spatial uncertainty. We thus propose to rely on the δ -FWER, that is the probabili-
15 ty of making a false discovery at a distance greater than δ from any true positive.
16 With this target measure in mind, we study the properties of ensembled clustered
17 inference algorithms which combine three techniques: spatially constrained cluster-
18 ing, statistical inference, and ensembling to aggregate several clustered inference
19 solutions. We show that ensembled clustered inference algorithms control the δ -
20 FWER under standard assumptions for δ equal to the largest cluster diameter. We
21 complement the theoretical analysis with empirical results, demonstrating accurate
22 δ -FWER control and decent power achieved by such inference algorithms.

23 Keywords: Clustering; High-dimension; Linear model; Spatial tolerance; Statistical in-
24 ference; Structured data; Support recovery.

1 Introduction

High-dimensional setting. High-dimensional regression corresponds to a setting where the number of covariates (or features) p exceeds the number of samples n . It notably occurs when searching for conditional associations among some high-dimensional observations and some outcome of interest: the *target*. Typical examples of the high-dimensional setting include inference problems on high-resolution images, where one aims at pixel- or voxel-level analysis, *e.g.*, in neuroimaging [Norman et al., 2006, De Martino et al., 2008], astronomy [Richards et al., 2009], but also in other fields where covariates display a spatial structure *e.g.*, in genomics [Balding, 2006, Dehman et al., 2015]. In all these examples, it actually turns out that not only $n < p$ but even $n \ll p$ and the covariates are spatially structured because of the physics of the problem or the measurements process. Because such high-dimensional data lead to high-variance results, probing statistical significance is important to give a level of confidence in the reported association. For this reason, the present analysis departs from traditional sparse modeling methods such as the Lasso [Tibshirani, 1996], that simply aim at selecting a good set of predictive covariates without considering statistical significance. In this context, a first approach is to consider the multivariate linear model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\varepsilon} ,$$

where the target is denoted by $\mathbf{y} \in \mathbb{R}^n$, the design matrix by $\mathbf{X} \in \mathbb{R}^{n \times p}$, the parameter vector by $\boldsymbol{\beta}^* \in \mathbb{R}^p$ and the random error vector by $\boldsymbol{\varepsilon} \in \mathbb{R}^n$. The aim is to infer $\boldsymbol{\beta}^*$, with statistical guarantees on the estimate, in particular regarding the support, *i.e.*, the set of covariates with non-zero importance.

Statistical inference on individual parameters. In high-dimensional settings, standard statistical inference methodology does not apply, but numerous methods have recently been proposed to recover the non-zero parameters of $\boldsymbol{\beta}^*$ with statistical guarantees. Many methods rely on resampling: bootstrap procedures [Bach, 2008, Chatterjee and Lahiri, 2011, Liu and Yu, 2013], perturbation resampling-based procedures [Minier et al., 2011], stability selection procedures [Meinshausen and Bühlmann, 2010] and randomized sample splitting [Wasserman and Roeder, 2009, Meinshausen et al., 2009]. All of these approaches suffer from limited power. Contrarily to the screening/inference procedure, post-selection inference procedures generally merge the screening and inference steps into one and then use all the samples [Berk et al., 2013, Lockhart et al., 2014, Lee et al., 2016, Tibshirani et al., 2016], resulting in potentially more powerful tests than sample splitting. Yet, these approaches do not scale well with large p . Another family of methods rely on debiasing procedures: the most prominent examples are corrected ridge [Bühlmann, 2013] and desparsified Lasso [Zhang and Zhang, 2014, van de Geer et al., 2014, Javanmard and Montanari, 2014] which is an active area of research [Javanmard and Montanari, 2018, Bellec and Zhang, 2019, Celentano et al., 2020]. Additionally, knockoff filters [Barber and Candès, 2015, Candès et al., 2018] consist in creating noisy “fake” copies of the original variables, and checking which original variables are selected

48 prior to the fake ones. Finally, a general framework for statistical inference in sparse
49 high-dimensional models has been proposed recently [Ning and Liu, 2017].

50 **Failure of existing statistical inference methods.** In practice, in the $n \ll p$ setting
51 we consider, the previous methods are not well adapted as they are often powerless or
52 computationally intractable. In particular, the number of predictive parameters (*i.e.*, the
53 support size) denoted $s(\beta^*)$ can be greater than the number of samples even in the
54 sparse setting, where $s(\beta^*) \ll p$. There is an underlying identifiability problem: in
55 general, one cannot retrieve all predictive parameters, as highlighted *e.g.*, in Wainwright
56 [2009]. Beyond the fact that statistical inference is impossible when $p \gg n$, the problem
57 is aggravated by the following three effects. First, as outlined above, dense covariate
58 sampling leads to high values for p and induces high correlation among covariates, further
59 challenging the conditions for recovery, as shown in Wainwright [2009]. Second, when
60 testing for several multiple hypothesis, the correction cost is heavy [Dunn, 1961, Westfall
61 and Young, 1993, Benjamini and Hochberg, 1995]; for example with Bonferroni correction
62 [Dunn, 1961], p-values are corrected by a factor p when testing every covariate. This
63 make this type of inference methods powerless in our settings (see Fig. 3 for instance).
64 Third, the above approaches are at least quadratic or cubic in the support size, hence
65 become prohibitive whenever both p and n are large.

66 **Combining clustering and inference.** Nevertheless, in these settings, variables of-
67 ten reflect some underlying spatial structure, such as smoothness. For example, in med-
68 ical imaging, an image has a 3D structure and a given voxel is highly correlated with
69 neighboring voxels; in genomics, there exist blocks of Single Nucleotide Polymorphisms
70 (SNPs) that tend to be jointly predictive or not. Hence, β^* can in general be assumed
71 to share the same structure: among several highly correlated covariates, asserting that
72 only one is important to predict the target seems meaningless, if not misleading.

73 A computationally attractive solution that alleviates high dimensionality is to group
74 correlated neighboring covariates. This step can be understood as a design compression:
75 it produces a closely related, yet reduced version of the original problem (see *e.g.*, Park
76 et al. [2006], Varoquaux et al. [2012], Hoyos-Idrobo et al. [2018]). Inference combined
77 with a fixed clustering has been proposed by Bühlmann et al. [2013] and can overcome
78 the dimensionality issue, yet this study does not provide procedures that derive cluster-
79 wise confidence intervals or p-values. Moreover, in most cases groups (or clusters) are not
80 pre-determined nor easily identifiable from data, and their estimation simply represents
81 a local optimum among a huge, non-convex space of solutions. It is thus problematic to
82 base inference upon such an arbitrary data representation. Inspired by this dimension
83 reduction approach, we have proposed [Chevalier et al., 2018] the ensemble of clustered
84 desparsified Lasso (EnCluDL) procedure that exhibits strong empirical performances
85 [Chevalier et al., 2021] in terms of support recovery even when $p \gg n$. EnCluDL is
86 an ensembled clustered inference algorithm, *i.e.*, it combines a spatially constrained
87 clustering procedure that reduces the problem dimension, an inference procedure that
88 performs statistical inference at the cluster level, and an ensembling method that ag-

gregates several cluster-level solutions. Concerning the inference step, the desparsified Lasso [Zhang and Zhang, 2014, van de Geer et al., 2014, Javanmard and Montanari, 2014] was preferred over other high-dimensional statistical inference procedures based on the comparative study of Dezeure et al. [2015] and on the research activity around it [Dezeure et al., 2017, Javanmard and Montanari, 2018, Bellec and Zhang, 2019, Cellentano et al., 2020]; however, it is possible to use another inference procedure that produces a p-value family controlling the classical FWER. By contrast, we did not consider the popular knockoff procedure [Barber and Candès, 2015, Candès et al., 2018], that does not produce p-values and does not control the family-wise error rate (FWER). However, an extension of the knockoffs to FWER-type control was proposed by Janson and Su [2016]. It does not control the standard FWER but another relaxed version of the FWER called k -FWER. As it is a relevant alternative to ensembled clustered inference algorithms, we have included it in our empirical comparison (see Section 5). In Nguyen et al. [2020], a variant of the knockoffs is proposed to control the FWER, but it does not handle large- p problems. Another extension that produces p-value, called conditional randomization test, has been presented in Candès et al. [2018], but its computational cost is prohibitive. Additionally, Meinshausen [2015] provides “group-bound” confidence intervals, corresponding to confidence intervals on the ℓ_1 -norm of several parameters, without making further assumptions on the design matrix. However, this method is known to be conservative in practice [Mitra and Zhang, 2016, Javanmard and Montanari, 2018]. Finally, hierarchical testing [Mandozzi and Bühlmann, 2016, Blanchard and Geman, 2005, Meinshausen, 2008] also leverages this clustering/inference combination but in a different way. Their approach consists in performing significance tests along the tree of a hierarchical clustering algorithm starting from the root node and descending subsequently into children of rejected nodes. This procedure has the drawback of being constrained by the clustering tree, which is often not available, thus replaced by some noisy estimate.

Contributions. Producing a cluster-wise inference is not completely satisfactory as it relies on an arbitrary clustering choice. Instead, we look for methods that derive covariate-wise statistics enabling support identification with a spatially relaxed false detection control. In that regard, our first contribution is to present a generalization of the FWER called δ -FWER, that takes into account a spatial tolerance of magnitude δ for the false discoveries. Then, our main contribution is to prove that ensembled clustered inference algorithms control the δ -FWER under reasonable assumptions for a given tolerance parameter δ . Finally, we apply the ensembled clustered inference scheme to the desparsified Lasso leading to the EnCluDL algorithm and conduct an empirical study: we show that EnCluDL exhibits a good statistical power in comparison with alternative procedures and we verify that it displays the expected δ -FWER control.

Notation. Throughout the remainder of this article, for any $p \in \mathbb{N}^*$, we write $[p]$ for the set $\{1, \dots, p\}$. For a vector β , β_j refers to its j -th coordinate. For a matrix \mathbf{X} , $\mathbf{X}_{i,\cdot}$ refers to the i -th row and $\mathbf{X}_{\cdot,j}$ to the j -th column and $\mathbf{X}_{i,j}$ refers to the element in the

130 i -th row and j -th column.

131 2 Model and data assumptions

132 2.1 Generative models of high-dimensional data: random fields

133 In the setting that we consider, we assume that the covariates come with a natural rep-
 134 resentation in a discretized metric space, generally the discretized 2D or 3D Euclidean
 135 space. In such settings, discrete random fields are convenient to model the random vari-
 136 ables representing the covariates. Indeed, denoting by $\mathbf{X} = (\mathbf{X}_{i,j})_{i \in [n], j \in [p]}$ the random
 137 design matrix, where n is the number of samples and p the number of covariates, the
 138 rows $(\mathbf{X}_{i,\cdot})_{i \in [n]}$ are sampled from a random field defined on a discrete domain.

139 2.2 Gaussian random design model and high dimensional settings

We assume that the covariates are independent and identically distributed and follow
 a centered Gaussian distribution, *i.e.*, for all $i \in [n]$, $\mathbf{X}_{i,\cdot} \sim \mathcal{N}(0_p, \mathbf{\Sigma})$ where $\mathbf{\Sigma}$ is the
 covariance matrix of the covariates. Our aim is to derive confidence bounds or p-values
 on the coefficients of the parameter vector denoted by β^* , under the Gaussian linear
 model:

$$\mathbf{y} = \mathbf{X}\beta^* + \varepsilon, \quad (1)$$

140 where $\mathbf{y} \in \mathbb{R}^n$ is the target, $\mathbf{X} \in \mathbb{R}^{n \times p}$ is the (random) design matrix, $\beta^* \in \mathbb{R}^p$ is the
 141 vector or parameters, and $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2 \mathbf{I}_n)$ is the noise vector with standard deviation
 142 $\sigma_\varepsilon > 0$. We make the assumption that ε is independent of \mathbf{X} .

143 2.3 Data structure

144 Since the covariates have a natural representation in a metric space, we assume that the
 145 spatial distances between covariates are known. With a slight abuse of notation, the
 146 distance between covariates j and k is denoted by $d(j, k)$ for $(j, k) \in [p] \times [p]$ and the
 147 correlation between covariates j and k is given by $\text{Cor}(\mathbf{X}_{\cdot,j}, \mathbf{X}_{\cdot,k}) = \mathbf{\Sigma}_{j,k} / \sqrt{\mathbf{\Sigma}_{j,j} \mathbf{\Sigma}_{k,k}}$. We
 148 now introduce a key structural assumption: two covariates at a spatial distance smaller
 149 than δ are positively correlated.

150 **Assumption 2.1.** *The covariates verify the spatial homogeneity assumption with dis-*
 151 *tance parameter $\delta > 0$ if, for all $(j, k) \in [p] \times [p]$, $d(j, k) \leq \delta$ implies that $\mathbf{\Sigma}_{j,k} \geq 0$.*

152 Under model (1), each coordinate of the parameter vector β^* links one covariate to
 153 the target. Then, β^* has the same underlying organization as the covariates and is also
 154 called weight map in these settings. Defining its *support* as $S(\beta^*) = \{j \in [p] : \beta_j^* \neq 0\}$
 155 and its cardinal as $s(\beta^*) = |S(\beta^*)|$, we assume that the true model is sparse, meaning
 156 that β^* has a small number of non-zero entries, *i.e.*, $s(\beta^*) \ll p$. The complementary
 157 of $S(\beta^*)$ in $[p]$ is called the *null region* and is denoted by $N(\beta^*)$, *i.e.*, $N(\beta^*) = \{j \in$

158 $[p] : \beta_j^* = 0\}$. Additionally to the sparse assumption, we assume that β^* is (spatially)
 159 smooth. To reflect sparsity and smoothness, we introduce another key assumption:
 160 weights associated with close enough covariates share the same sign, zero being both
 161 positive and negative.

162 **Assumption 2.2.** *The weight vector β^* verifies the sparse-smooth assumption with*
 163 *distance parameter $\delta > 0$ if, for all $(j, k) \in [p] \times [p]$, $d(j, k) \leq \delta$ implies that $\text{sign}(\beta_j^*) =$*
 164 *$\text{sign}(\beta_k^*)$.*

165 Equivalently, the sparse-smooth assumption with parameter δ holds if the distance
 166 between the two closest weights of opposite sign is larger than δ . In Fig. 2-(a), we give
 167 an example of a weight map verifying the sparse-smooth assumption with $\delta = 2$.

168 3 Statistical control with spatial tolerance

169 Under the spatial assumption we have discussed, discoveries that are closer than δ from
 170 the true support are not considered as false discoveries: inference at a resolution finer
 171 than δ might be unrealistic. This means that δ can be interpreted as a tolerance param-
 172 eter on the (spatial) support we aim at recovering. Then, we introduce a new metric
 173 closely related to the FWER that takes into account spatial tolerance and we call it δ -
 174 family wise error rate (δ -FWER). A similar extension of the false discovery rate (FDR)
 175 has been introduced by Cheng et al. [2020], Nguyen et al. [2019], Gimenez and Zou [2019],
 176 but, to the best of our knowledge, this has not been considered yet for the FWER. In
 177 the following, we consider a general estimator $\hat{\beta}$ that comes with p-values, testing the
 178 nullity of the corresponding parameters, denoted by $\hat{p} = (\hat{p}_j)_{j \in [p]}$. Also, we denote by
 179 $S(\hat{\beta}) \subset [p]$ a general estimate of the support $S(\beta^*)$ derived from the estimator $\hat{\beta}$.

Definition 3.1 (δ -null hypothesis). *For all $j \in [p]$, the δ -null hypothesis for the j -th*
covariates, $H_0^\delta(j)$, states that all other covariates at distance less than δ have a zero
weight in the true model (1); the alternative hypothesis is denoted $H_1^\delta(j)$:

$$H_0^\delta(j) : \text{“for all } k \in [p] \text{ such that } d(j, k) \leq \delta, \beta_k^* = 0 \text{” ,}$$

$$H_1^\delta(j) : \text{“there exists } k \in [p] \text{ such that } d(j, k) \leq \delta \text{ and } \beta_k^* \neq 0 \text{” .}$$

180 Thus, we say that a δ -type 1 error is made if a null covariate $j \in [p]$ is selected,
 181 i.e., $j \in S(\hat{\beta})$, while $H_0^\delta(j)$ holds true. Taking $\delta = 0$ recovers the usual null-hypothesis
 182 $H_0(j) : \beta_j^* = 0$ and usual type 1 error.

Definition 3.2 (Control of the δ -type 1 error). *The p-value related to the j -th covariate*
denoted by \hat{p}_j controls the δ -type 1 error if, under $H_0^\delta(j)$, for all $\alpha \in (0, 1)$, we have:

$$\mathbb{P}(\hat{p}_j \leq \alpha) \leq \alpha ,$$

183 where \mathbb{P} is the probability distribution with respect to the random dataset of observations
 184 (\mathbf{X}, \mathbf{y}) .

Definition 3.3 (δ -null region). *The set of indexes of covariates verifying the δ -null hypothesis is called the δ -null region and is denoted by $N^\delta(\beta^*)$ (or simply N^δ):*

$$N^\delta(\beta^*) = \{j \in [p] : \text{for all } k \in [p], d(j, k) \leq \delta \text{ implies that } \beta_k^* = 0\} .$$

185 When $\delta = 0$ the δ -null region is simply the null region : $N^0(\beta^*) = N(\beta^*)$. We also
 186 point out the nested property of δ -null regions with respect to δ : for $0 \leq \delta_1 \leq \delta_2$ we
 187 have $N^{\delta_2}(\beta^*) \subseteq N^{\delta_1}(\beta^*) \subseteq N(\beta^*)$ (see Fig. 2-(d) for an example of δ -null region).

Definition 3.4 (Rejection region). *Given a family of p -values $\hat{p} = (\hat{p}_j)_{j \in [p]}$ and a threshold $\alpha \in (0, 1)$, the rejection region, $R_\alpha(\hat{p})$, is the set of indexes having a p -value lower than α :*

$$R_\alpha(\hat{p}) = \{j \in [p] : \hat{p}_j \leq \alpha\} .$$

Definition 3.5 (δ -type 1 error region). *Given a family of p -values $\hat{p} = (\hat{p}_j)_{j \in [p]}$ and a threshold $\alpha \in (0, 1)$, the δ -type 1 error region at level α is $\mathcal{E}_\alpha^\delta$, the set of indexes belonging both to the δ -null region and to the rejection region at level α . We also refer to this region as the erroneous rejection region at level α with tolerance δ :*

$$\mathcal{E}_\alpha^\delta(\hat{p}) = N^\delta \cap R_\alpha(\hat{p}) .$$

188 When $\delta = 0$ the δ -type 1 error region recovers the type 1 error region which is
 189 denoted by $\mathcal{E}_\alpha(\hat{p})$. Again, one can verify a nested property: for $0 \leq \delta_1 \leq \delta_2$ we have
 190 $\mathcal{E}_\alpha^{\delta_2}(\hat{p}) \subseteq \mathcal{E}_\alpha^{\delta_1}(\hat{p}) \subseteq \mathcal{E}_\alpha(\hat{p})$.

Definition 3.6 (δ -family wise error rate). *Given a family of p -values $\hat{p} = (\hat{p}_j)_{j \in [p]}$ and a threshold $\alpha \in (0, 1)$, the δ -FWER at level α with respect to the family \hat{p} , denoted $\delta\text{-FWER}_\alpha(\hat{p})$, is the probability that the δ -type 1 error region at level α is not empty:*

$$\delta\text{-FWER}_\alpha(\hat{p}) = \mathbb{P}(|\mathcal{E}_\alpha^\delta(\hat{p})| \geq 1) = \mathbb{P}(\min_{j \in N^\delta} \hat{p}_j \leq \alpha) .$$

Definition 3.7 (δ -FWER control). *We say that the family of p -values $\hat{p} = (\hat{p}_j)_{j \in [p]}$ controls the δ -FWER if, for all $\alpha \in (0, 1)$:*

$$\delta\text{-FWER}_\alpha(\hat{p}) \leq \alpha .$$

191 When $\delta = 0$ the δ -FWER is the usual FWER. Additionally, for $0 \leq \delta_1 \leq \delta_2$, one can
 192 verify that $\delta_2\text{-FWER}_\alpha(\hat{p}) \leq \delta_1\text{-FWER}_\alpha(\hat{p}) \leq \text{FWER}_\alpha(\hat{p})$. Thus, δ -FWER control is a
 193 weaker property than usual FWER control.

194 4 δ -FWER control with clustered inference algorithms

195 4.1 Clustered inference algorithms

196 A clustered inference algorithm consists in partitioning the covariates into groups (or
 197 clusters) before applying a statistical inference procedure. In Sec. 4.1, we describe a

198 standard clustered inference algorithm that produces a (corrected) p-value family on the
 199 parameters of the model (1). In this algorithm, in addition to the observations (\mathbf{X}, \mathbf{y}) , we
 200 take as input the transformation matrix $\mathbf{A} \in \mathbb{R}^{p \times C}$ which maps and averages covariates
 201 into C clusters. The `statistical_inference` function corresponds to a given statistical
 202 inference procedure that takes as inputs the clustered data \mathbf{Z} and the target \mathbf{y} and
 203 produces valid p-values for every cluster. If $C < n$, least squares are suitable, otherwise,
 204 procedures such as multi-sample split [Wasserman and Roeder, 2009, Meinshausen et al.,
 205 2009], corrected ridge [Bühlmann, 2013] or desparsified Lasso [Zhang and Zhang, 2014,
 206 van de Geer et al., 2014, Javanmard and Montanari, 2014] might be relevant whenever
 207 their assumptions are verified. Then, the computed p-values are corrected for multiple
 208 testing by multiplying by a factor C . Finally, covariate-wise p-values are inherited from
 209 the corresponding cluster-wise p-values.

Algorithm 1 Clustered inference

```

input :  $\mathbf{X} \in \mathbb{R}^{n \times p}, \mathbf{y} \in \mathbb{R}^n, \mathbf{A} \in \mathbb{R}^{p \times C}$ 
 $\mathbf{Z} = \mathbf{X}\mathbf{A}$  // compressed design matrix
 $\hat{p}^{\mathcal{G}} = \text{statistical\_inference}(\mathbf{Z}, \mathbf{y})$  // uncorrected cluster-wise p-values
 $\hat{q}^{\mathcal{G}} = C \times \hat{p}^{\mathcal{G}}$  // corrected cluster-wise p-values
for  $j = 1, \dots, p$  do
  |  $\hat{q}_j = \hat{q}_c^{\mathcal{G}}$  if  $j$  in cluster  $c$  // corrected covariate-wise p-values
return  $\hat{q} = (\hat{q}_j)_{j \in [p]}$  // family of corrected covariate-wise p-values

```

Algorithm 2 Ensembled clustered inference

```

input :  $\mathbf{X} \in \mathbb{R}^{n \times p}, \mathbf{y} \in \mathbb{R}^n$ 
param :  $C, B$ 
for  $b = 1, \dots, B$  do
  |  $\mathbf{X}^{(b)} = \text{sampling}(\mathbf{X})$  // sampling rows of  $\mathbf{X}$ 
  |  $\mathbf{A}^{(b)} = \text{clustering}(q, \mathbf{X}^{(b)})$  // transformation matrix
  |  $\hat{q}^{(b)} = \text{clustered\_inference}(\mathbf{X}, \mathbf{y}, \mathbf{A}^{(b)})$  // families of corr. covariate-wise p-val.
for  $j = 1, \dots, p$  do
  |  $\hat{q}_j = \text{ensembling}(\{\hat{q}_j^{(b)}, b \in [B]\})$  // aggregated corrected covariate-wise p-values
return  $\hat{q} = (\hat{q}_j)_{j \in [p]}$  // family of aggregated corrected covariate-wise p-values

```

210 Ensembled clustered inference algorithms correspond to the ensembling of several
 211 clustered inference solutions for different choice of clusterings using the p-value aggre-
 212 gation proposed by Meinshausen et al. [2009]. In Sec. 4.1, we give a standard ensembled
 213 clustered inference algorithm that produces a (corrected) p-value family on the param-
 214 eters of the model (1). In this algorithm, the `sampling` function corresponds to a subsam-
 215 pling of the data, *i.e.*, a subsampling of the rows of \mathbf{X} . The `clustering` function derives
 216 a choice of clustering in C clusters, it produces a transformation matrix $\mathbf{A}^{(b)} \in \mathbb{R}^{p \times C}$

217 that should vary for each bootstrap $b \in [B]$ since the subsampled data $\mathbf{X}^{(b)}$ varies. Once
 218 the clustering inference steps are completed, the `ensembling` function aggregates the B
 219 (corrected) p-value families into a single one.

220 **Fig. 1** can help the reader to better understand the organization of the next sections,
 221 aiming eventually at establishing the δ -FWER control property of the clustered inference
 222 and ensembled clustered inference algorithms.

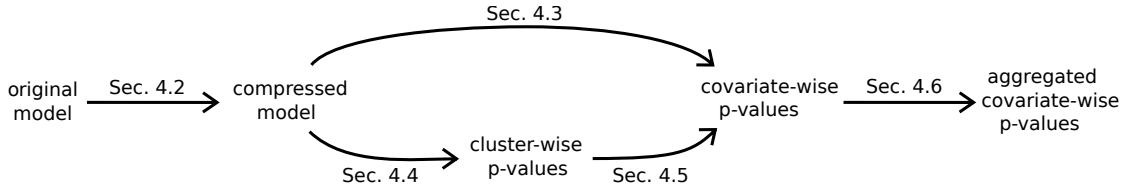


Figure 1: Organization of **Section 4**.

223 4.2 Compressed representation

The motivation for using groups of covariates that are spatially concentrated is to reduce the dimension while preserving large-scale data structure. The number of groups is denoted by $C < p$ and, for $r \in [q]$, we denote by G_r the r -th group. The collection of all the groups is denoted by $\mathcal{G} = \{G_1, G_2, \dots, G_C\}$ and forms a partition of $[p]$. Every group representative variable is defined by the average of the covariates it contains. Then, denoting by $\mathbf{Z} \in \mathbb{R}^{n \times C}$ the compressed random design matrix that contains the group representative variables in columns and, without loss of generality, assuming a suitable ordering of the columns of \mathbf{X} , dimension reduction can be written:

$$\mathbf{Z} = \mathbf{X}\mathbf{A} \quad , \quad (2)$$

where $\mathbf{A} \in \mathbb{R}^{p \times q}$ is the transformation matrix defined by:

$$\mathbf{A} = \begin{bmatrix} \alpha_1 - \alpha_1 & 0 - 0 & \dots & 0 - 0 \\ 0 - 0 & \alpha_2 - \alpha_2 & \dots & 0 - 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 - 0 & 0 - 0 & \dots & \alpha_C - \alpha_C \end{bmatrix} \quad ,$$

where $\alpha_c = 1/|G_c|$ for all $c \in [C]$. Consequently, the distribution of the i -th row of \mathbf{Z} is given by $\mathbf{Z}_{i,\cdot} \sim \mathcal{N}_q(0, \mathbf{\Upsilon})$, where $\mathbf{\Upsilon} = \mathbf{A}^\top \mathbf{\Sigma} \mathbf{A}$. The correlation between the groups $r \in [q]$ and $l \in [q]$ is given by $\text{Cor}(\mathbf{Z}_{\cdot,r}, \mathbf{Z}_{\cdot,l}) = \Upsilon_{r,l} / \sqrt{\Upsilon_{r,r} \Upsilon_{l,l}}$. As mentioned in [Bühlmann et al. \[2013\]](#), because of the Gaussian assumption in (1), we have the following compressed representation:

$$\mathbf{y} = \mathbf{Z}\boldsymbol{\theta}^* + \boldsymbol{\eta} \quad , \quad (3)$$

224 where $\boldsymbol{\theta}^* \in \mathbb{R}^q$, $\boldsymbol{\eta} \sim \mathcal{N}(0, \sigma_\eta^2 \mathbf{I}_n)$, $\sigma_\eta \geq \sigma_\varepsilon > 0$ and $\boldsymbol{\eta}$ is independent of \mathbf{Z} .

225 **Remark 4.1.** *Dimension reduction is not the unique desirable effect of clustering with*
 226 *regards to statistical inference. Indeed, this clustering-based design compression also*
 227 *generally improves the conditioning of the problem. Assumptions needed for valid statis-*
 228 *tical inference are thus more likely to be met. For more details about this conditioning*
 229 *enhancement, the reader may refer to [Bühlmann et al. \[2013\]](#).*

230 4.3 Properties of the compressed model weights

231 We now give a property of the weights of the compressed problem which is a consequence
 232 of [Bühlmann et al. \[2013, Proposition 4.3\]](#).

Proposition 4.1. *Considering the Gaussian linear model in (1) and assuming:*

- (i) for all $c \in [C]$, for all $(j, k) \in (G_c)^2$, $\Sigma_{j,k} \geq 0$,
- (ii) for all $c \in [C]$, for all $c' \in [C] \setminus \{c\}$, $\Upsilon_{c,c'} = 0$,
- (iii) for all $c \in [C]$, $(\beta_j^* \geq 0$ for all $j \in G_c)$ or $(\beta_j^* \leq 0$ for all $j \in G_c)$,

233 then, in the compressed representation (3), for $c \in [C]$, $\theta_c^* \neq 0$ if and only if there exists
 234 $j \in G_c$ such that $\beta_j^* \neq 0$. If such an index j exists then $\text{sign}(\theta_r^*) = \text{sign}(\beta_j^*)$.

235 *Proof.* See Supplement [E.1](#). □

236 Assumption (i) states that the covariates in a group are all positively correlated. Let
 237 us define the group diameter (or cluster diameter) of G_c by the distance that separates
 238 its two most distant covariates, *i.e.*, $\text{Diam}(G_c) = \max\{d(j, k) : (j, k) \in (G_c)^2\}$ and the
 239 clustering diameter of \mathcal{G} by the largest group diameter, *i.e.*, $\text{Diam}(\mathcal{G}) = \max\{\text{Diam}(G_c) :$
 240 $c \in [C]\}$. In [Fig. 2-\(b\)](#), we propose a clustering of the initial weight map in [Fig. 2-\(a\)](#) for
 241 which the clustering diameter is equal to 2 for the ℓ_1 distance. Assumption (i) notably
 242 holds when $\text{Diam}(\mathcal{G}) \leq \delta$ under the spatial homogeneity assumption ([Ass. 2.1](#)) with pa-
 243 rameter δ . Assumption (ii) assumes independence of the groups. A sufficient condition
 244 is when the covariates covariance matrix Σ is block diagonal, with blocks coinciding
 245 with the group structure; *i.e.*, assumption (ii) holds when covariates of different groups
 246 are independent. In practice, this assumption may be unmet, and we relax it in Sup-
 247 plement [B](#). Assumption (iii) states that all the weights in a group share the same sign.
 248 This is notably the case when the clustering diameter is smaller than δ and the weight
 249 map satisfies the sparse-smooth assumption ([Ass. 2.2](#)) with parameter δ . For instance,
 250 a clustering-based compressed representation of the weight map in [Fig. 2-\(a\)](#) is given in
 251 [Fig. 2-\(c\)](#).

252 4.4 Statistical inference on the compressed model

To perform the statistical inference on the compressed problem (3), we could consider
 any statistical inference procedure that produces cluster-wise p-values $\hat{p}^{\mathcal{G}} = (\hat{p}_c^{\mathcal{G}})_{c \in [C]}$,
 given a choice of clustering \mathcal{G} , that control the type 1 error. More precisely, for any
 $c \in [C]$, under $H_0(G_c)$, *i.e.*, the null hypothesis which states that θ_c^* is equal to zero

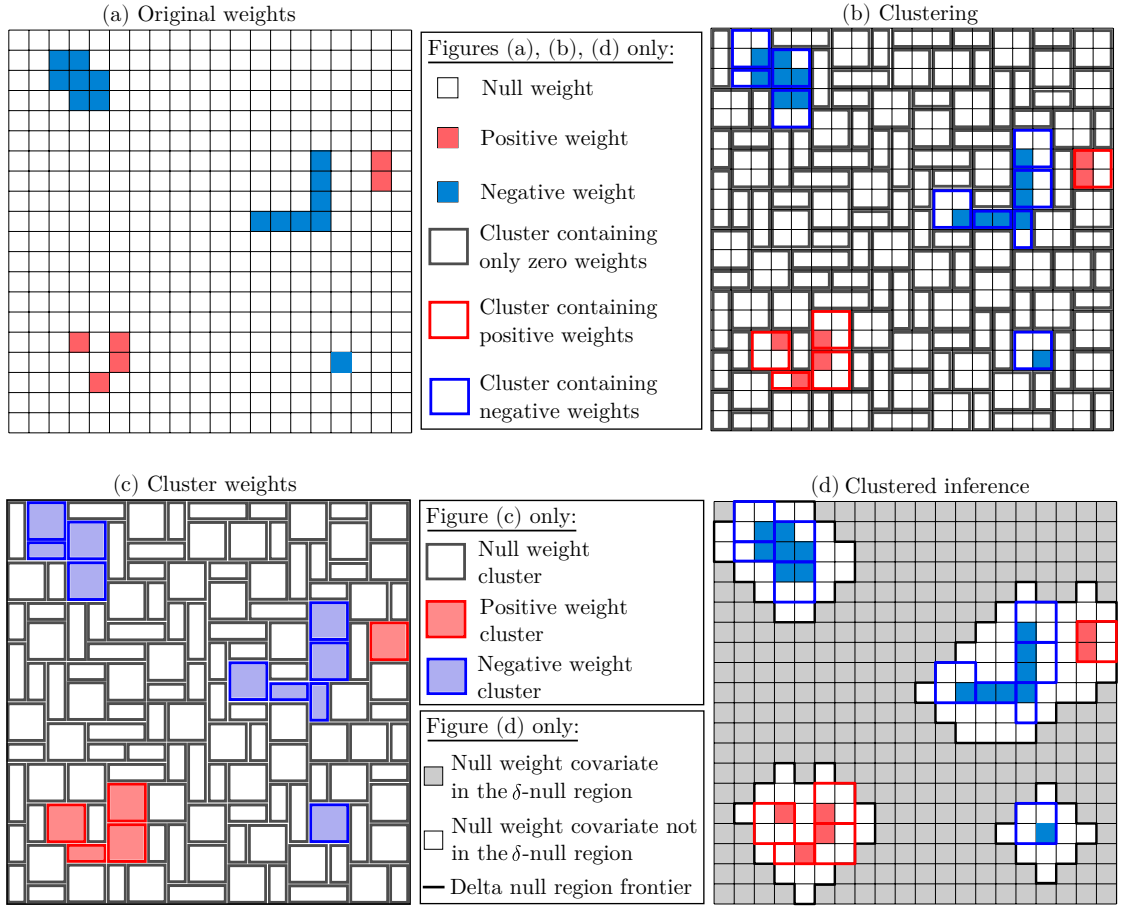


Figure 2: Clustered inference mechanism on 2D-spatially structured data. Item a: Example of weight map with a 2D-structure. Voxels represent covariates, with blue (resp. red) corresponding to negative (resp. positive) weights; others are null weights. Item b: Arbitrary choice of spatially constrained clustering with a diameter of $\delta = 2$ units for the ℓ_1 -distance. Rectangles delimited by black lines represent clusters that contain only zero-weight covariates. Blue (resp. red) rectangles refer to clusters that contain negative-weight (resp. positive) covariates. Item c: Compressed model weights: under the assumptions of [Prop. 4.1](#), the cluster weights share the same signs as the covariate weights they contain. Blue (resp. red) rectangles correspond to negative-weight (resp. positive-weights) clusters. Item d: The grey area corresponds to the δ -null region ($\delta = 2$). Under the same assumptions, the non-zero weight groups have no intersection with the δ -null region.

in the compressed model, we assume that the p-value associated with the c -th cluster verifies:

$$\mathbb{P}(\hat{p}_c^{\mathcal{G}} \leq \alpha) \leq \alpha . \quad (4)$$

To correct for multiple comparisons, we consider Bonferroni correction [Dunn, 1961] which is a conservative procedure but has the advantage of being valid without any additional assumptions. Furthermore, here the correction factor is only equal to the number of groups, not the number of covariates. Then, the family of corrected cluster-wise p-values $\hat{q}^{\mathcal{G}} = (\hat{q}_c^{\mathcal{G}})_{c \in [C]}$ is defined by:

$$\hat{q}_c^{\mathcal{G}} = \min\{1, C \times \hat{p}_c^{\mathcal{G}}\} . \quad (5)$$

Let us denote by $N_{\mathcal{G}}(\boldsymbol{\theta}^*)$ (or simply $N_{\mathcal{G}}$) the null region in the compressed problem for a given choice of clustering \mathcal{G} , *i.e.*, $N_{\mathcal{G}}(\boldsymbol{\theta}^*) = \{c \in [C] : \boldsymbol{\theta}_c^* = 0\}$. Then, for all $\alpha \in (0, 1)$:

$$\text{FWER}_{\alpha}(\hat{q}^{\mathcal{G}}) = \mathbb{P}(\min_{c \in N_{\mathcal{G}}} \hat{q}_c^{\mathcal{G}} \leq \alpha) \leq \alpha . \quad (6)$$

253 This means that the cluster-wise p-value family $\hat{q}^{\mathcal{G}}$ controls FWER.

254 4.5 De-grouping

Given the families of cluster-wise p-values $\hat{p}^{\mathcal{G}}$ and corrected p-values $\hat{q}^{\mathcal{G}}$ as defined in (10) and (5), our next aim is to derive families of p-values and corrected p-values related to the covariates of the original problem. To construct these families, we simply set the (corrected) p-value of the j -th covariate to be equal to the (corrected) p-value of its corresponding group:

$$\begin{aligned} \text{for all } j \in [p], \quad \hat{p}_j &= \sum_{c \in [C]} \mathbb{1}_{\{j \in G_c\}} \hat{p}_c^{\mathcal{G}} , \\ \text{for all } j \in [p], \quad \hat{q}_j &= \sum_{c \in [C]} \mathbb{1}_{\{j \in G_c\}} \hat{q}_c^{\mathcal{G}} . \end{aligned} \quad (7)$$

255 **Proposition 4.2.** *Under the assumptions of Prop. 4.1 and assuming that the clustering*
256 *diameter is smaller than δ , then:*

(i) *elements of the family \hat{p} defined in (7) control the δ -type 1 error:*

$$\text{for all } j \in N^{\delta}, \text{ for all } \alpha \in (0, 1), \mathbb{P}(\hat{p}_j \leq \alpha) \leq \alpha ,$$

(ii) *the family \hat{q} defined in (7) controls the δ -FWER:*

$$\text{for all } \alpha \in (0, 1), \mathbb{P}(\min_{j \in N^{\delta}} \hat{q}_j \leq \alpha) \leq \alpha .$$

257 *Proof.* See Supplement E.2. □

258 The previous de-grouping properties can be seen in Fig. 2-(d). Roughly, since all
259 the clusters that intersect the δ -null region have low p-value with low probability, one
260 can conclude that all the covariates of the δ -null region also have low p-value with low
261 probability.

262 **4.6 Ensembling**

263 Let us consider B families of corrected p-values that control the δ -FWER. For any
 264 $b \in [B]$, we denote by $\hat{q}^{(b)}$ the b -th family of corrected p-values. Then, we show that
 265 the ensembling method proposed in Meinshausen et al. [2009] yields a family that also
 266 enforces δ -FWER control.

Proposition 4.3. *Assume that, for $b \in [B]$, the p-value families $\hat{q}^{(b)}$ control the δ -FWER. Then, for any $\gamma \in (0, 1)$, the ensembled p-value family $\tilde{q}(\gamma)$ defined by:*

$$\text{for all } j \in [p], \tilde{q}_j(\gamma) = \min \left\{ 1, \gamma\text{-quantile} \left(\left\{ \frac{\hat{q}_j^{(b)}}{\gamma} : b \in [B] \right\} \right) \right\}, \quad (8)$$

267 controls the δ -FWER.

268 *Proof.* See Supplement E.3. □

269 **4.7 δ -FWER control**

270 We can now state our main result: the clustered inference and ensembled clustered
 271 inference algorithms control the δ -FWER.

272 **Theorem 4.1.** *Assume the model given in (1) and that the data structure assumptions,
 273 Ass. 2.1 and Ass. 2.2, are satisfied for a distance parameter larger than δ . Assume
 274 that all the clusterings considered have a diameter smaller than δ . Assume that the
 275 uncorrelated cluster assumption, i.e., assumption (ii) of Prop. 4.1, is verified for each
 276 clustering and further assume that the statistical inference performed on the compressed
 277 model (3) is valid, i.e., (4) holds. Then, the p-value family obtained from the clustered
 278 inference algorithm controls the δ -FWER. Additionally, the p-value family derived by the
 279 ensembled clustered inference algorithm controls the δ -FWER.*

280 *Proof.* See Supplement E.4. □

281 **Remark 4.2.** *When the type 1 error control offered by the statistical inference proce-
 282 dure is only asymptotic, the result stated by Theorem 4.1 remains true asymptotically.
 283 This is notably the case when using desparsified Lasso: under the assumptions of Theo-
 284 rem 4.1 and the assumptions specific to desparsified Lasso (cf. Supplement A), ensemble
 285 of clustered desparsified Lasso (EnCluDL) controls the δ -FWER asymptotically.*

286 **5 Numerical Simulations**

287 **5.1 CluDL and EnCluDL**

288 For testing the (ensembled) clustered inference algorithms, we have decided to make
 289 the inference step using the desparsified Lasso [Zhang and Zhang, 2014, van de Geer
 290 et al., 2014, Javanmard and Montanari, 2014] leading to the clustered desparsified Lasso

291 (CluDL) and the ensemble of clustered desparsified Lasso (EnCluDL) algorithms that
292 were first presented in [Chevalier et al. \[2018\]](#).

293 In Supplement [A](#), we detail the assumptions and refinements that occur when choos-
294 ing the desparsified Lasso to perform the statistical inference step. A notable difference
295 is the fact that all the results becomes asymptotic. In Supplement [C](#), we present a
296 diagram illustrating the mechanism of EnCluDL and analyse its numerical complexity.

297 5.2 2D Simulation

298 We run a series of simulations on 2D data in order to give empirical evidence of the
299 theoretical properties of CluDL and EnCluDL and compare their recovery properties
300 with two other procedures. For an easier visualization of the results, we consider one
301 central scenario, whose parameters are written in **bold** in the following of this section,
302 with several variations, changing only one parameter at a time.

303 In all these simulations, the feature space considered is a 2D square with edge length
304 $H = 40$ leading to $p = H^2 = 1\,600$ covariates, with a sample size $n \in \{50, \mathbf{100}, 200, 400\}$.
305 To construct β^* , we define a 2D weight map $\tilde{\beta}^*$ with four active regions (as illustrated in
306 [Fig. 3](#)) and then flatten $\tilde{\beta}^*$ to a vector β^* of size p . Each active region is a square of width
307 $h \in \{2, \mathbf{4}, 6, 8\}$, leading to a size of support of 1%, **4%**, 9% or 16%. To construct the
308 design matrix, we first build a 2D data matrix $\tilde{\mathbf{X}}$ by drawing p random normal vectors
309 of size n that are spatially smoothed with a 2D Gaussian filter to create a correlation
310 structure related to the covariates' spatial organization. The same flattening process as
311 before is used to get the design matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$. The intensity of the spatial smoothing
312 is adjusted to achieve a correlation between two adjacent covariates (local correlation)
313 of $\rho \in \{0.5, \mathbf{0.75}, 0.9, 0.95\}$. We also set the noise standard deviation $\sigma_\varepsilon \in \{1, \mathbf{2}, 3, 4\}$,
314 which corresponds to a signal to noise ratio (SNR) $\text{SNR}_y \in \{6.5, \mathbf{3.5}, 2.2, 1.5\}$, where the
315 SNR is defined by $\text{SNR}_y = \|\mathbf{X}\beta^*\|_2 / \|\varepsilon\|_2$. For each scenario, we run 100 simulations
316 to derive meaningful statistics. A Python implementation of the simulations and proce-
317 dures presented in this paper is available on <https://github.com/ja-che/hidimstat>.
318 Regarding the clustering step in CluDL and EnCluDL, we used a spatially constrained
319 agglomerative clustering algorithm with Ward criterion. This algorithm is popular in
320 many applications [[Varoquaux et al., 2012](#), [Dehman et al., 2015](#)], as it tends to create
321 compact, balanced clusters. Since the optimal number of clusters C is unknown a priori,
322 we have tested several values $C \in [100; 400]$. A smaller C generally improves recovery,
323 but entails a higher spatial tolerance. Following theoretical considerations, we compute
324 the largest cluster diameter for every value of C and set δ to this value. We obtained the
325 couples $(C, \delta) \in \{(100, 8), (200, 6), (300, 5), (400, 4)\}$. The tolerance region is represented
326 in [Fig. 3](#) for $\delta = 6$. Concerning EnCluDL, we took a number of bootstraps B equal to
327 25 as we observed that it was sufficient to benefit from most of the effect of clustering
328 randomization.

329 **5.3 Alternative methods**

330 We compare the recovering properties of CluDL and EnCluDL with two other proce-
 331 dures: desparsified Lasso and knockoffs. Contrarily to CluDL and EnCluDL, none of
 332 these includes a compression step. The version of the desparsified Lasso we have tested
 333 is the one presented in van de Geer et al. [2014], that outputs p-values. Using Bonferroni
 334 correction it controls the classical FWER at any desired rate. The original version of
 335 knockoffs [Barber and Candès, 2015, Candès et al., 2018] only controls the false discov-
 336 ery rate (FDR) which is a weaker control than the classical FWER. Yet Janson and Su
 337 [2016] modifies the covariate selection process leading to a procedure that controls the
 338 k -FWER, *i.e.*, the probability of making at least k false discoveries. We tested this last
 339 extension of knockoffs. Depending on the nominal rate at which we want to control the
 340 k -FWER, the choice of k is not arbitrary. More precisely, if we want a k -FWER control
 341 at 10%, we need to tolerate $k = 4$ at least, otherwise the estimated support would always
 342 be empty.

343 Since k -FWER and δ -FWER controls are both weaker than the usual FWER control
 344 whenever $k > 1$ and $\delta > 0$, one can expect desparsified Lasso to be less powerful than
 345 knockoffs, CluDL and EnCluDL. Besides, there is no relation between k -FWER and
 346 δ -FWER controls when $k > 1$ and $\delta > 0$, hence it is not possible to establish which one
 347 is less prohibitive for support recovery. However, when data are spatially structured,
 348 δ -FWER control might be more relevant since it controls the very undesirable far-from-
 349 support false discoveries.

350 **5.4 Results**

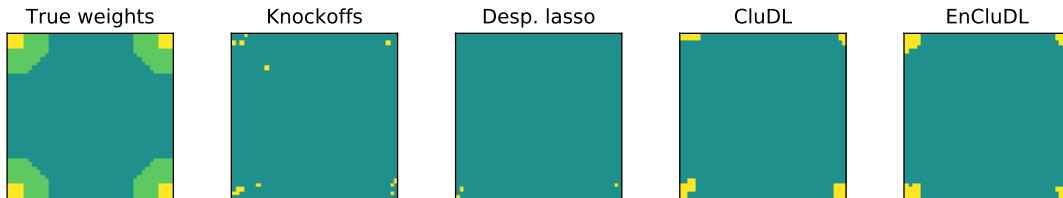


Figure 3: True support and estimated support for the first seed of the central scenario. Left: The support in yellow is composed of four regions of width $h = 4$ covariates. The tolerance region in green surrounds the support, its width is $\delta = 6$ covariates. The remaining covariates in blue form the δ -null region. Others: The yellow squares are the covariates selected by each method. Knockoffs selects few covariates when controlling the k -FWER at 10% for $k = 4$. Desparsified Lasso only retrieves 3 covariates when controlling the FWER at 10%. For $C = 200$, CluDL and EnCluDL have good power and control the δ -FWER at 10% for $\delta = 6$.

351 In Fig. 3, we plot the maps estimated by knockoffs, desparsified Lasso, CluDL and
 352 EnCluDL for $C = 200$ when solving the first seed of the central scenario simulation.
 353 Regarding knockoffs and desparsified Lasso solutions, we notice that the power is low

354 and the methods select few covariates in each predictive region. The CluDL method is
 355 more powerful and recovers groups of covariates that correspond more closely to the true
 356 weights. However, the shape of the CluDL solution depends on the clustering choice.
 357 The EnCluDL solution seems even more powerful than the CluDL one and recovers
 358 groups of covariates that correspond almost perfectly to the true weights. Both CluDL
 359 and EnCluDL are only accurate up to the spatial tolerance which is $\delta = 6$, but EnCluDL
 360 fits the ground truth more tightly.

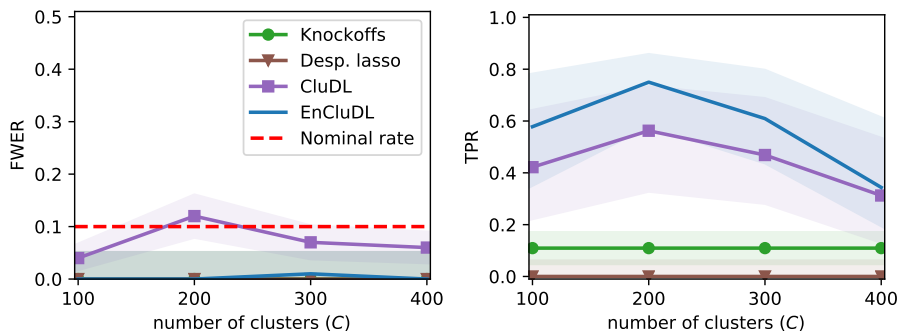


Figure 4: Results for fixed simulation parameters corresponding to the central scenario simulation. The green line with circles correspond to knockoffs, the brown line with triangles is the desparsified Lasso, the purple squared line correspond to CluDL and the blue plain line is EnCluDL. Left: Empirical FWER for desparsified Lasso, k -FWER for knockoffs and δ -FWER for CluDL and EnCluDL. The 80% confidence intervals are obtained by Binomial approximation. Right: Median true positive rate (TPR) for all the procedures, together with 80% confidence interval obtained by taking the first decile and last decile TPR.

361 In Fig. 4, we focus on the central scenario to get more insight about the statistical
 362 properties of the methods and the influence of the C hyper-parameter for CluDL and
 363 EnCluDL. First, we observe that all methods reach the targeted control: desparsified
 364 Lasso controls the FWER, knockoffs control the k -FWER and, CluDL and EnCluDL
 365 control the δ -FWER. Second, considering the true positive rates (TPR), we notice that
 366 the methods that do not integrate a compression step, *i.e.*, knockoffs and desparsified
 367 Lasso, have a limited statistical power due to $n \ll p$. However, CluDL has decent power
 368 and EnCluDL improves over CluDL thanks to clustering randomization. Finally, CluDL
 369 and EnCluDL are flexible with respect to the choice of C since the TPR varies quite
 370 slowly with C .

371 We have also studied the influence of the simulation parameters by varying one pa-
 372 rameter of the central scenario. The corresponding results are available in Supplement D.
 373 The main conclusion gained from these complementary results is the fact that, up to the
 374 limit given by the desired spatial tolerance δ , the choice of C should be made in function
 375 of the data structure. More precisely, good clustering creates clusters that are weakly
 376 correlated and contains covariates that are highly correlated. This observation is linked
 377 to assumption (ii) of Prop. 4.1.

378 6 Discussion

379 When $n \ll p$, statistical inference on predictive model parameters is a hard problem.
380 However, when the data are spatially structured, we have shown that ensembled clustered
381 inference procedures are attractive, as they exhibit statistical guarantees and good power.
382 The price to pay is to accept that inference is only accurate up to spatial distance δ
383 corresponding to the clustering diameter, thus replacing FWER with δ -FWER control
384 guarantees.

385 One of the most obvious field of application of this class of algorithms is neuroscience
386 where it can be used to solve source localization problems. In that regards, a wide
387 empirical validation of EnCluDL has been conducted in [Chevalier et al. \[2021\]](#) including
388 fMRI data experiments. Also, an extension of EnCluDL was proposed in [Chevalier](#)
389 [et al. \[2020\]](#) to address the magneto/electroencephalography source localization problem
390 which involves spatio-temporal data.

391 With EnCluDL, the statistical inference step is performed by the desparsified Lasso.
392 In [Nguyen et al. \[2019\]](#), another ensembled clustered inference method that leverages
393 the knockoff technique [[Barber and Candès, 2015](#)] leading to a procedure called ECKO
394 has been tested. However, formal δ -FDR control guarantees have not been established
395 yet for this model. It would be also quite natural to try other inference techniques such
396 as the (distilled) conditional randomization test [[Candès et al., 2018](#), [Liu and Janson,](#)
397 [2020](#)].

398 In the present work, we have only considered the linear regression setup. However,
399 combining the same algorithmic scheme with statistical inference solutions for gener-
400 alized linear models, we could extend this work to the logistic regression setup. This
401 would extend the usability of ensembled clustered inference to many more application
402 settings.

403 Acknowledgement

404 This study has been funded by Labex DigiCosme (ANR-11-LABEX-0045-DIGICOSME)
405 as part of the program "Investissement d'Avenir" (ANR-11-IDEX-0003-02), by the Fast-
406 Big project (ANR-17-CE23-0011) and the KARAIB AI Chair (ANR-20-CHIA-0025-01).
407 This study has also been supported by the European Union's Horizon 2020 research and
408 innovation program (Grant Agreement No. 945539, Human Brain Project SGA3).

409 Supplementary material

410 Supplementary material available online includes an analysis of the technical assump-
411 tions and refinements that occur when choosing the desparsified Lasso to perform the
412 statistical inference step in Supplement [A](#), a diagram summarizing EnCluDL and a study
413 of the complexity of EnCluDL in Supplement [C](#), a proposition for relaxing assumption
414 (ii) of [Prop. 4.1](#) in Supplement [B](#), complementary results for studying the influence of
415 the simulation parameters in Supplement [D](#) and the proofs in Supplement [E](#).

416 References

- 417 F. R. Bach. Bolasso: model consistent lasso estimation through the bootstrap. In
418 *Proceedings of the 25th international conference on Machine learning*, pages 33–40,
419 2008. 2
- 420 D. J. Balding. A tutorial on statistical methods for population association studies.
421 *Nature reviews genetics*, 7(10):781–791, 2006. 2
- 422 R. F. Barber and E. Candès. Controlling the false discovery rate via knockoffs. *Ann.*
423 *Statist.*, 43(5):2055–2085, 10 2015. 2, 4, 15, 17
- 424 P. C. Bellec and C.-H. Zhang. De-biasing the lasso with degrees-of-freedom adjustment.
425 *arXiv preprint arXiv:1902.08885*, 2019. 2, 4
- 426 Y. Benjamini and Y. Hochberg. Controlling the False Discovery Rate: A Practical and
427 Powerful Approach to Multiple Testing. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 57
428 (1):289–300, 1995. 3
- 429 R. Berk, L. Brown, A. Buja, K. Zhang, and L. Zhao. Valid post-selection inference. *Ann.*
430 *Statist.*, 41(2):802–837, 2013. 2
- 431 G. Blanchard and D. Geman. Hierarchical testing designs for pattern recognition. *The*
432 *Annals of Statistics*, 33(3):1155–1202, 2005. 4
- 433 P. Bühlmann. Statistical significance in high-dimensional linear models. *Bernoulli*, 19
434 (4):1212–1242, 09 2013. 2, 8
- 435 P. Bühlmann, P. Rütimann, S. van de Geer, and C.-H. Zhang. Correlated variables
436 in regression: Clustering and sparse estimation. *Journal of Statistical Planning and*
437 *Inference*, 143(11):1835 – 1858, 2013. 3, 9, 10, 24, 28, 29
- 438 E. Candès, Y. Fan, L. Janson, and J. Lv. Panning for gold: ‘model-X’ knockoffs for high
439 dimensional controlled variable selection. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 80
440 (3):551–577, 2018. 2, 4, 15, 17
- 441 M. Celentano, A. Montanari, and Y. Wei. The lasso with general gaussian designs with
442 applications to hypothesis testing. *arXiv preprint arXiv:2007.13716*, 2020. 2, 4
- 443 A. Chatterjee and S. N. Lahiri. Bootstrapping lasso estimators. *J. Amer. Statist. Assoc.*,
444 106(494):608–625, 2011. 2
- 445 D. Cheng, Z. He, A. Schwartzman, et al. Multiple testing of local extrema for detection
446 of change points. *Electron. J. Stat.*, 14(2):3705–3729, 2020. 6
- 447 J.-A. Chevalier, J. Salmon, and B. Thirion. Statistical inference with ensemble of clus-
448 tered desparsified lasso. In *International Conference on Medical Image Computing*
449 *and Computer-Assisted Intervention*, pages 638–646. Springer, 2018. 3, 14

- 450 J.-A. Chevalier, A. Gramfort, J. Salmon, and B. Thirion. Statistical control for spatio-
451 temporal meg/eeg source imaging with desparsified multi-task lasso. In *Thirty-fourth*
452 *Conference on Neural Information Processing Systems*, 2020. 17
- 453 J.-A. Chevalier, T.-B. Nguyen, J. Salmon, G. Varoquaux, and B. Thirion. Decoding with
454 confidence: Statistical control on decoder maps. *NeuroImage*, page 117921, 2021. 3,
455 17
- 456 F. De Martino, G. Valente, N. Staeren, J. Ashburner, R. Goebel, and E. Formisano.
457 Combining multivariate voxel selection and support vector machines for mapping and
458 classification of fMRI spatial patterns. *Neuroimage*, 43(1):44–58, 2008. 2
- 459 A. Dehman, C. Ambroise, and P. Neuvial. Performance of a blockwise approach in
460 variable selection using linkage disequilibrium information. *BMC bioinformatics*, 16
461 (1):148, 2015. 2, 14
- 462 R. Dezeure, P. Bühlmann, L. Meier, and N. Meinshausen. High-dimensional inference:
463 Confidence intervals, p -values and R-Software hdi. *Statist. Sci.*, 30(4):533–558, 2015.
464 4
- 465 R. Dezeure, P. Bühlmann, and C.-H. Zhang. High-dimensional simultaneous inference
466 with the bootstrap. *Test*, 26(4):685–719, 2017. 4
- 467 O. J. Dunn. Multiple comparisons among means. *J. Amer. Statist. Assoc.*, 56(293):
468 52–64, 1961. 3, 12
- 469 J. R. Gimenez and J. Zou. Discovering conditionally salient features with statistical
470 guarantees. *International Conference on Machine Learning*, pages 2290–2298, 2019. 6
- 471 A. Hoyos-Idrobo, G. Varoquaux, J. Kahn, and B. Thirion. Recursive nearest agglomera-
472 tion (rena): fast clustering for approximation of structured signals. *IEEE transactions*
473 *on pattern analysis and machine intelligence*, 41(3):669–681, 2018. 3
- 474 L. Janson and W. Su. Familywise error rate control via knockoffs. *Electron. J. Stat.*, 10
475 (1):960–975, 2016. 4, 15
- 476 A. Javanmard and A. Montanari. Confidence intervals and hypothesis testing for high-
477 dimensional regression. *J. Mach. Learn. Res.*, 15:2869–2909, 2014. 2, 4, 8, 13, 22
- 478 A. Javanmard and A. Montanari. Debiasing the lasso: Optimal sample size for Gaussian
479 designs. *Ann. Statist.*, 46(6A):2593–2622, 2018. 2, 4
- 480 J. Lee, D. Sun, Y. Sun, and J. Taylor. Exact post-selection inference, with application
481 to the lasso. *Ann. Statist.*, 44(3):907–927, 2016. 2
- 482 H. Liu and B. Yu. Asymptotic properties of lasso+ mls and lasso+ ridge in sparse
483 high-dimensional linear regression. *Electron. J. Stat.*, 7:3124–3169, 2013. 2

- 484 M. Liu and L. Janson. Fast and powerful conditional randomization testing via distilla-
485 tion. *arXiv preprint arXiv:2006.03980*, 2020. 17
- 486 R. Lockhart, J. Taylor, R. J. Tibshirani, and R. Tibshirani. A significance test for the
487 lasso. *Ann. Statist.*, 42(2):413, 2014. 2
- 488 J. Mandozzi and P. Bühlmann. Hierarchical testing in the high-dimensional setting with
489 correlated variables. *J. Amer. Statist. Assoc.*, 111(513):331–343, 2016. 4
- 490 N. Meinshausen. Hierarchical testing of variable importance. *Biometrika*, 95(2):265–278,
491 2008. 4
- 492 N. Meinshausen. Group bound: confidence intervals for groups of variables in sparse
493 high dimensional regression without assumptions on the design. *J. R. Stat. Soc. Ser.*
494 *B Stat. Methodol.*, pages 923–945, 2015. 4
- 495 N. Meinshausen and P. Bühlmann. Stability selection. *J. R. Stat. Soc. Ser. B Stat.*
496 *Methodol.*, 72:417–473, 2010. 2
- 497 N. Meinshausen, L. Meier, and P. Bühlmann. P-values for high-dimensional regression.
498 *J. Amer. Statist. Assoc.*, 104(488):1671–1681, 2009. 2, 8, 13, 30
- 499 J. Minnier, L. Tian, and T. Cai. A perturbation method for inference on regularized
500 regression estimates. *J. Amer. Statist. Assoc.*, 106(496):1371–1382, 2011. 2
- 501 R. Mitra and C.-H. Zhang. The benefit of group sparsity in group inference with de-
502 biased scaled group lasso. *Electron. J. Stat.*, 10(2):1829–1873, 2016. 4
- 503 E. Ndiaye, O. Fercoq, A. Gramfort, V. Leclère, and J. Salmon. Efficient smoothed
504 concomitant lasso estimation for high dimensional regression. In *Journal of Physics:*
505 *Conference Series*, volume 904, page 012006. IOP Publishing, 2017. 22
- 506 T.-B. Nguyen, J.-A. Chevalier, and B. Thirion. Ecko: Ensemble of clustered knock-
507 offs for robust multivariate inference on fMRI data. In *International Conference on*
508 *Information Processing in Medical Imaging*, pages 454–466. Springer, 2019. 6, 17
- 509 T.-B. Nguyen, J.-A. Chevalier, B. Thirion, and S. Arlot. Aggregation of multiple knock-
510 offs. In *International Conference on Machine Learning*, pages 7283–7293. PMLR,
511 2020. 4
- 512 Y. Ning and H. Liu. A general theory of hypothesis tests and confidence regions for
513 sparse high dimensional models. *Ann. Statist.*, 45(1):158–195, 2017. 3
- 514 K. A. Norman, S. M. Polyn, G. J. Detre, and J. V. Haxby. Beyond mind-reading: multi-
515 voxel pattern analysis of fMRI data. *Trends in cognitive sciences*, 10(9):424–430, 2006.
516 2
- 517 M. Y. Park, T. Hastie, and R. Tibshirani. Averaged gene expressions for regression.
518 *Biostatistics*, 8(2):212–227, 05 2006. 3

- 519 S. Reid, R. Tibshirani, and J. Friedman. A study of error variance estimation in lasso
520 regression. *Statistica Sinica*, pages 35–67, 2016. [22](#)
- 521 J.W. Richards, P.E. Freeman, A.B. Lee, and C.M. Schafer. Exploiting low-dimensional
522 structure in astronomical spectra. *The Astrophysical Journal*, 691(1):32, 2009. [2](#)
- 523 R. Tibshirani. Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. Ser. B*
524 *Stat. Methodol.*, 58(1):267–288, 1996. [2](#)
- 525 R. J. Tibshirani, J. Taylor, R. Lockhart, and R. Tibshirani. Exact post-selection infer-
526 ence for sequential regression procedures. *J. Amer. Statist. Assoc.*, 111(514):600–620,
527 2016. [2](#)
- 528 S. van de Geer, P. Bühlmann, Y. Ritov, and R. Dezeure. On asymptotically optimal
529 confidence regions and tests for high-dimensional models. *Ann. Statist.*, 42(3):1166–
530 1202, 2014. [2](#), [4](#), [8](#), [13](#), [15](#), [22](#)
- 531 G. Varoquaux, A. Gramfort, and B. Thirion. Small-sample brain mapping: sparse recov-
532 ery on spatially correlated designs with randomization and clustering. In *International*
533 *Conference on Machine Learning*, 2012. [3](#), [14](#)
- 534 M. J. Wainwright. Sharp thresholds for high-dimensional and noisy sparsity recovery
535 using ℓ_1 -constrained quadratic programming (lasso). *IEEE Trans. Image Process.*, 55
536 (5):2183–2202, 2009. [3](#)
- 537 L. Wasserman and K. Roeder. High-dimensional variable selection. *Ann. Statist.*, 37
538 (5A):2178–2201, 2009. [2](#), [8](#)
- 539 P. H. Westfall and S. S. Young. *Resampling-based multiple testing: Examples and meth-*
540 *ods for p-value adjustment*, volume 279. John Wiley & Sons, 1993. [3](#)
- 541 G. Yu and J. Bien. Estimating the error variance in a high-dimensional linear model.
542 *Biometrika*, 106(3):533–546, 2019. [22](#)
- 543 C.-H. Zhang and S. S. Zhang. Confidence intervals for low dimensional parameters in
544 high dimensional linear models. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 76(1):217–242,
545 2014. [2](#), [4](#), [8](#), [13](#), [22](#)

546
547

Supplementary material for “Spatially relaxed inference on high-dimensional linear models”

548

A Desparsified Lasso on the compressed model

549
550
551
552
553
554
555
556
557
558
559

Here, we clarify the assumptions and refinements that occur when choosing the desparsified Lasso as the procedure that performs the statistical inference on the compressed model. The desparsified Lasso was first developed in Zhang and Zhang [2014] and Javanmard and Montanari [2014], and thoroughly analyzed in van de Geer et al. [2014]. Following notation in Eq. (3), the true support in the compressed model is denoted by $S(\boldsymbol{\theta}^*) = \{c \in [C] : \boldsymbol{\theta}_c^* \neq 0\}$ and its cardinality by $s(\boldsymbol{\theta}^*) = |S(\boldsymbol{\theta}^*)|$. We also denote by $\boldsymbol{\Omega} \in \mathbb{R}^{C \times C}$ the inverse of the population covariance matrix of the groups, *i.e.*, $\boldsymbol{\Omega} = \boldsymbol{\Upsilon}^{-1}$. Then, for $c \in [C]$, the sparsity of the c -th row of $\boldsymbol{\Omega}$ (or c -th column) is $s(\boldsymbol{\Omega}_{c,\cdot}) = |S(\boldsymbol{\Omega}_{c,\cdot})|$, where $S(\boldsymbol{\Omega}_{c,\cdot}) = \{c' \in [C] : \boldsymbol{\Omega}_{c,c'} \neq 0\}$. We also denote the smallest eigenvalue of $\boldsymbol{\Upsilon}$ by $\phi_{\min}(\boldsymbol{\Upsilon}) > 0$. We can now state the assumptions required for probabilistic inference with desparsified Lasso [van de Geer et al., 2014]:

Theorem A.1 (Theorem 2.2 of van de Geer et al. [2014]). *Considering the model in Eq. (3) and assuming:*

- (i) $1/\phi_{\min}(\boldsymbol{\Upsilon}) = \mathcal{O}(1)$,
- (ii) $\max_{c \in [C]} \boldsymbol{\Upsilon}_{c,c} = \mathcal{O}(1)$,
- (iii) $s(\boldsymbol{\theta}^*) = o(\sqrt{n}/\log(C))$,
- (iv) $\max_{c \in [C]} (s(\boldsymbol{\Omega}_{c,\cdot})) = o(n/\log(C))$,

then, denoting by $\hat{\boldsymbol{\theta}}$ the desparsified Lasso estimator derived from the inference procedure described in van de Geer et al. [2014], the following holds:

$$\begin{aligned} \sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) &= \boldsymbol{\xi} + \boldsymbol{\zeta} , \\ \boldsymbol{\xi} | \mathbf{Z} &\sim \mathcal{N}(0_C, \sigma_\eta^2 \hat{\boldsymbol{\Omega}}) , \\ \|\boldsymbol{\zeta}\|_\infty &= o_{\mathbb{P}}(1) , \end{aligned}$$

560

where $\hat{\boldsymbol{\Omega}}$ is such that $\|\hat{\boldsymbol{\Omega}} - \boldsymbol{\Omega}\|_\infty = o_{\mathbb{P}}(1)$.

561
562
563
564

Remark A.1. *In Theorem A.1, to compute confidence intervals, the noise standard deviation σ_η in the compressed problem has to be estimated. We refer the reader to the surveys that are dedicated to this subject such as Reid et al. [2016], Ndiaye et al. [2017], Yu and Bien [2019].*

As argued in van de Geer et al. [2014], from Theorem A.1 we obtain asymptotic confidence intervals for the r -th element of $\boldsymbol{\theta}^*$ from the following equations, for all

$z_1 \in \mathbb{R}$ and $z_2 \in \mathbb{R}^+$:

$$\begin{aligned} \mathbb{P} \left[\frac{\sqrt{n}(\hat{\boldsymbol{\theta}}_c - \boldsymbol{\theta}_c^*)}{\sigma_\eta \sqrt{\hat{\boldsymbol{\Omega}}_{c,c}}} \leq z_1 \mid \mathbf{Z} \right] - \Phi(z_1) &= o_{\mathbb{P}}(1) , \\ \mathbb{P} \left[\frac{\sqrt{n}|\hat{\boldsymbol{\theta}}_c - \boldsymbol{\theta}_c^*|}{\sigma_\eta \sqrt{\hat{\boldsymbol{\Omega}}_{c,c}}} \leq z_2 \mid \mathbf{Z} \right] - (2\Phi(z_2) - 1) &= o_{\mathbb{P}}(1) , \end{aligned} \quad (9)$$

where $\Phi(\cdot)$ is the cumulative distribution function of the standard normal distribution. Thus, for each $c \in [C]$ one can provide a p-value that assesses whether or not $\boldsymbol{\theta}_c^*$ is equal to zero. In the case of a two-sided single test, for each $c \in [C]$, the p-value denoted by $\hat{p}_c^{\mathcal{G}}$ is:

$$\hat{p}_c^{\mathcal{G}} = 2 \left(1 - \Phi \left(\frac{\sqrt{n}|\hat{\boldsymbol{\theta}}_c|}{\sigma_\eta \sqrt{\hat{\boldsymbol{\Omega}}_{c,c}}} \right) \right) . \quad (10)$$

Under $H_0(G_c)$, from (9), we have, for any $\alpha \in (0, 1)$:

$$\begin{aligned} \mathbb{P}(\hat{p}_c^{\mathcal{G}} \leq \alpha \mid \mathbf{Z}) &= 1 - \mathbb{P} \left[\frac{\sqrt{n}|\hat{\boldsymbol{\theta}}_c|}{\sigma_\eta \sqrt{\hat{\boldsymbol{\Omega}}_{c,c}}} \leq \Phi^{-1} \left(1 - \frac{\alpha}{2} \right) \mid \mathbf{Z} \right] \\ &= \alpha + o_{\mathbb{P}}(1) . \end{aligned} \quad (11)$$

Then, (11) shows that the p-values $\hat{p}_c^{\mathcal{G}}$ asymptotically control type 1 errors. Using the Bonferroni correction, the family of corrected p-values $\hat{q}^{\mathcal{G}} = (\hat{q}_c^{\mathcal{G}})_{c \in [C]}$ remains defined by:

$$\hat{q}_c^{\mathcal{G}} = \min\{1, C \times \hat{p}_c^{\mathcal{G}}\} . \quad (12)$$

Then, for all $\alpha \in (0, 1)$:

$$\text{FWER}_\alpha(\hat{q}^{\mathcal{G}}) = \mathbb{P}(\min_{c \in N_G} \hat{q}_c^{\mathcal{G}} \leq \alpha \mid \mathbf{Z}) \leq \alpha + o_{\mathbb{P}}(1) . \quad (13)$$

565 Then, (13) shows that the p-value family $\hat{q}^{\mathcal{G}}$ asymptotically control FWER. Finally,
 566 we have shown that desparsified Lasso applied to a compressed version of the original
 567 problem provides cluster-wise p-value families $\hat{p}^{\mathcal{G}}$ and $\hat{q}^{\mathcal{G}}$ that control respectively the
 568 type 1 error and the FWER in the compressed model only asymptotically.

569 B Relaxing the uncorrelated clusters assumption

570 As noted in Sec. 4.3, assumption (ii) of Prop. 4.1 is often unmet in practice. Here, taking
 571 the particular case in which the inference step is performed by desparsified Lasso, we
 572 relax the assumption and show that it is still possible to compute an adjusted corrected p-
 573 value that asymptotically controls the δ -FWER. Hopefully, the technique used to derive

574 this relaxation would also be applicable to other parametric statistical inference methods
575 such as corrected ridge. To better understand the development made in this section,
576 the adjusted p-values of this section should be compared with the original p-values of
577 Supplement A. Note that, this extension is easy to integrate in the proof of the main
578 results Theorem 4.1 as it just requires to use the adjusted corrected p-value instead
579 of the original corrected p-value. Also, it does not provide much more insight about
580 clustered inference algorithms. This is why we have decided to keep this extension for
581 Supplementary Materials.

582 First, we replace Prop. 4.1 by the next proposition that is a consequence of Bühlmann
583 et al. [2013, Proposition 4.4].

584 **Proposition B.1.** *Considering the Gaussian linear model in (1) and assuming:*

585 (i) for all $c \in [C]$, for all $j, k \in G_c^2$, $\text{Cov}(\mathbf{X}_{\cdot,j}, \mathbf{X}_{\cdot,k} \mid \{\mathbf{Z}_{\cdot,c'} : c' \neq c\}) \geq 0$,

(ii.a) for all $c \in [C]$, there exists $\nu_c \in \mathbb{R}^+$ s.t. for all $j \in G_c$, for all $k \notin G_c$,

$$|\text{Cov}(\mathbf{X}_{\cdot,j}, \mathbf{X}_{\cdot,k} \mid \{\mathbf{Z}_{\cdot,c'} : c' \neq c\})| \leq \nu_c ,$$

586(ii.b) for all $c \in [C]$, there exists $\tau_c > 0$ s.t. $\text{Var}(\mathbf{Z}_{\cdot,c} \mid \{\mathbf{Z}_{\cdot,c'} : c' \neq c\}) \geq \tau_c$,

587(iii) for all $c \in [C]$, $\left(\text{for all } j \in G_c, \beta_j^* \geq 0\right)$ or $\left(\text{for all } j \in G_c, \beta_j^* \leq 0\right)$,

then, in the compressed representation (3), θ^* admits the following decomposition:

$$\theta^* = \tilde{\theta} + \kappa , \tag{14}$$

588 where, for all $c \in [C]$, $|\kappa_c| \leq (\nu_c / \tau_c) \|\beta^*\|_1$ and $\tilde{\theta}_c \neq 0$ if and only if there exists $j \in G_c$
589 such that $\beta_j^* \neq 0$. If such an index j exists then $\text{sign}(\tilde{\theta}_c) = \text{sign}(\beta_j^*)$.

590 *Proof.* See Supplement E.1. □

591 The assumptions (i) and (ii) in Prop. 4.1 are replaced by (i), (ii.a) and (ii.b) in
592 Prop. B.1. More precisely, instead of assuming that the covariates inside a group are
593 positively correlated, we assume that they are positively correlated conditionally to all
594 other groups. Also, we relax the more questionable assumption of groups independence;
595 we assume instead that the conditional covariance of two covariates of different groups
596 is bounded above (ii.a) and that the conditional variance of the group representative
597 variable is non-zero (ii.b). In practice, except when group representative variables are
598 linearly dependent, we can always find values for which (ii.a) and (ii.b) are verified, but
599 we would like the upper bound of (ii.a) as low as possible and the lower bound of (ii.b)
600 as high as possible. Finally, assumption (iii) remains unchanged.

Then, as done in Supplement A, we can build $\hat{\theta}$. Under the same assumptions,
Theorem A.1 is still valid and $\hat{\theta}$ still verifies (9). However, here we want to estimate $\tilde{\theta}$,
not θ^* . Combining Theorem A.1 and Prop. B.1, we can see $\hat{\theta}$ as a biased estimator of

$\tilde{\theta}$. To take this bias into account, we need to adjust the definition of the p-values given by (10). Let us assume that, for a given $a \in \mathbb{R}^+$,

$$\max_{c \in [C]} \left(\frac{\nu_c}{\tau_c \sqrt{\hat{\Omega}_{c,c}}} \right) \leq \frac{a \sigma_\varepsilon}{\|\beta^*\|_1} . \quad (15)$$

And, for all $c \in [C]$, let us define the adjusted p-values:

$$\hat{p}_c^{\mathcal{G}} = 2 \left(1 - \Phi \left(\sqrt{n} \left[\frac{|\hat{\theta}_c|}{\sigma_\eta \sqrt{\hat{\Omega}_{c,c}}} - a \right]_+ \right) \right) . \quad (16)$$

601 Let us denote by $q_{1-\frac{\alpha}{2}} = \Phi^{-1}(1 - \frac{\alpha}{2})$ the $1 - \frac{\alpha}{2}$ quantile of the standard Gaussian
602 distribution. Then, under $H_0(G_c)$, the hypothesis which states that $\beta_j^* = 0$ for $j \in G_c$
603 implying that $\tilde{\theta}_c = 0$, we have, for any $\alpha \in (0, 1)$:

$$\begin{aligned} \mathbb{P}(\hat{p}_c^{\mathcal{G}} \leq \alpha \mid \mathbf{Z}) &= 1 - \mathbb{P} \left[\sqrt{n} \left[\frac{|\hat{\theta}_c|}{\sigma_\eta \sqrt{\hat{\Omega}_{c,c}}} - a \right]_+ \leq q_{1-\frac{\alpha}{2}} \mid \mathbf{Z} \right] \\ &\leq 1 - \mathbb{P} \left[\sqrt{n} \left[\frac{|\hat{\theta}_c|}{\sigma_\eta \sqrt{\hat{\Omega}_{c,c}}} - \frac{\nu_c \|\beta^*\|_1}{\sigma_\varepsilon \tau_c \sqrt{\hat{\Omega}_{c,c}}} \right]_+ \leq q_{1-\frac{\alpha}{2}} \mid \mathbf{Z} \right] \\ &\leq 1 - \mathbb{P} \left[\sqrt{n} \left[\frac{|\hat{\theta}_c| - |\kappa_c|}{\sigma_\eta \sqrt{\hat{\Omega}_{c,c}}} \right]_+ \leq q_{1-\frac{\alpha}{2}} \mid \mathbf{Z} \right] \\ &= 1 - \mathbb{P} \left[\sqrt{n} \left[\frac{|\hat{\theta}_c| - |\theta_c^*|}{\sigma_\eta \sqrt{\hat{\Omega}_{c,c}}} \right]_+ \leq q_{1-\frac{\alpha}{2}} \mid \mathbf{Z} \right] \\ &\leq 1 - \mathbb{P} \left[\sqrt{n} \frac{|\hat{\theta}_c - \theta_c^*|}{\sigma_\eta \sqrt{\hat{\Omega}_{c,c}}} \leq q_{1-\frac{\alpha}{2}} \mid \mathbf{Z} \right] \\ &= \alpha + o_{\mathbb{P}}(1) . \end{aligned} \quad (17)$$

604 Finally, we have built a cluster-wise adjusted p-value family that asymptotically exhibits,
605 with low probability ($< \alpha$), low value ($< \alpha$) for the clusters which contain only zero
606 weight covariates. To complete the proof in the case of correlated clusters, one can
607 proceed as in uncorrelated cluster case taking (16) instead of (10).

608 Now, let us come back to the interpretation and choice for the constant a . In
609 **Prop. B.1**, we have shown that, when groups are not independent, a group weight in the
610 compressed model can be non-zero even if the group only contains zero weight covariates.
611 However, the absolute value of the weight of such a group is necessarily upper bounded.
612 We thus introduce $a \in \mathbb{R}^+$ in (16) to increase the p-values by a relevant amount and

613 keep statistical guarantees concerning the non-discovery of a such group. The value of a
 614 depends on the physics of the problem and on the choice of clustering. While the physics
 615 of the problem is fixed, the choice of clustering has a strong impact on the left term of
 616 (15) and a "good" choice of clustering results in a lower a (less correction). To estimate
 617 a , we need to find an upper bound of $\|\beta^*\|_1$, a lower bound of σ_ε and to estimate the
 618 left term of (15). In practice, to compute p-values, we took $a = 0$ since the formula in
 619 (10) was already conservative for all the problems we considered.

620 C EnCluDL

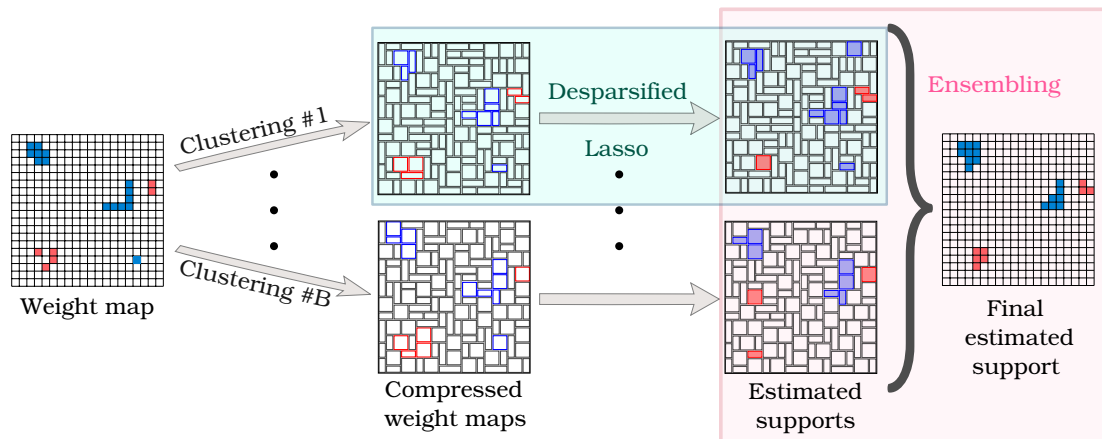


Figure 5: Summary of the mechanism of ensemble of clustered desparsified Lasso (EnCluDL). EnCluDL combines three algorithmic steps: a clustering procedure, the desparsified Lasso statistical inference procedure to derive p-value maps, and an ensembling method that synthesizes several p-value maps into one.

621 Computationally, to derive the EnCluDL solution we must solve B independent
 622 CluDL problems, making the global problem embarrassingly parallel; nevertheless, we
 623 could run the CluDL algorithm on standard desktop stations without parallelization
 624 with $n = 400$, $p \approx 10^5$, $C = 500$ and $B = 25$ in less than 10 minutes. Note that, the
 625 clustering step being much quicker than the inference step, p has a very limited impact
 626 on the total computation time.

627 The complexity for solving the Lasso depends significantly on the choice of solver,
 628 we then give the complexity in numbers of Lasso. The complexity for solving EnCluDL
 629 is given by the complexity of the resolution of $\mathcal{O}(B \times C)$ Lasso problems with n samples
 630 and C covariates, *i.e.*, with clustering. It is noteworthy that the complexity for solving
 631 the desparsified Lasso on the original problem is given by the complexity of the resolution
 632 of $\mathcal{O}(p)$ Lasso problems with n samples and p covariates, *i.e.*, without clustering. Then,
 633 EnCluDL should be much faster than the desparsified Lasso whenever $p \gg C$.

634 **D Complementary simulation results**

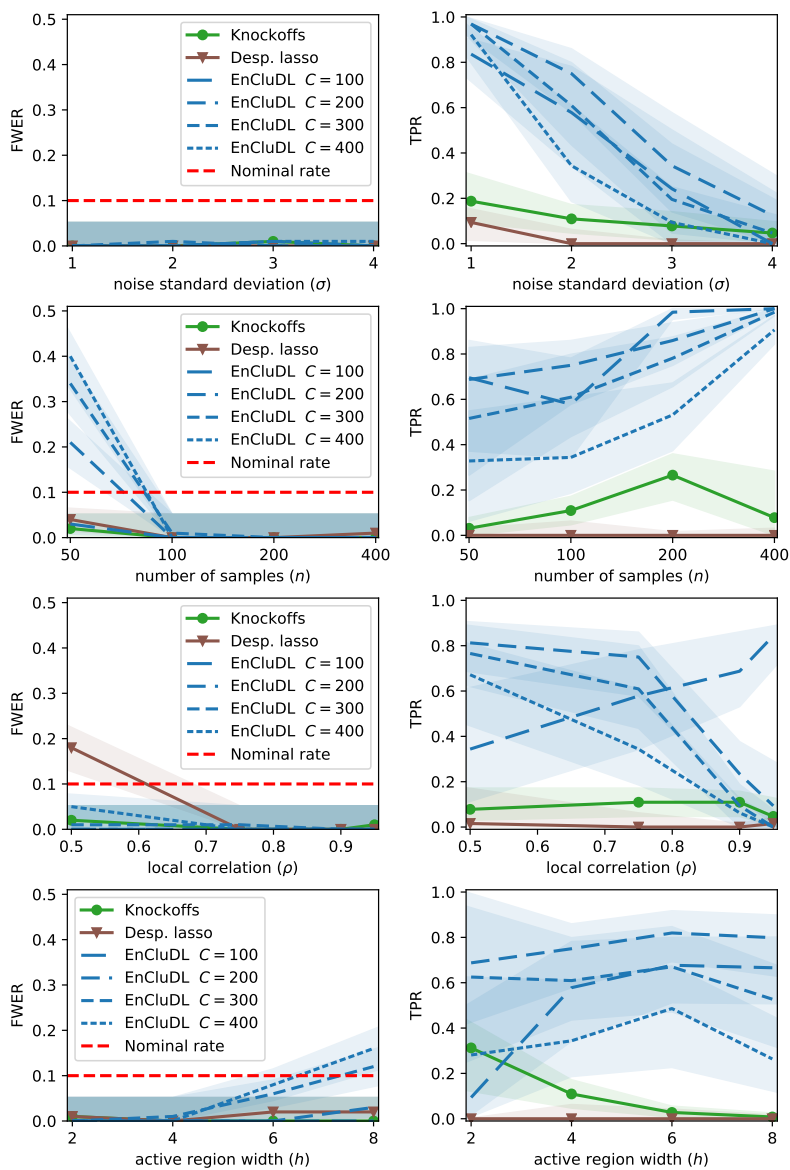


Figure 6: Results for various simulation parameters. The green line with circles correspond to the knockoffs, the brown line with triangle is the desparsified lasso, the dashed blue lines are for EnCluDL with length of the dashes increasing when C diminishes: large dashes are for $C = 100$, medium for $C = 200$, small for $C = 300$, tiny for $C = 400$. We compute the same FWER and TPR quantities as in Fig. 4, and the same 80% confidence intervals: by Binomial approximation for the FWER and taking first and last deciles for the TPR.

635 In Fig. 6, we study the influence of the simulation parameters by varying one param-
 636 eter of the central scenario at a time. We vary the noise standard deviation, the number
 637 of samples, the local correlation and the size of the support. For a better readability of
 638 the figures, we do not analyze the results of CLuDL since it is expected to be always
 639 a bit less powerful than EnCluDL while showing a similar behavior. First, we look at
 640 the plots where we vary the noise standard deviation σ . We observe that the methods
 641 reach the targeted FWER control and notice that EnCluDL benefits more strongly from
 642 the decrease of σ regarding support recovery. Second, we analyze the results for vari-
 643 ous sample sizes (n) values. Concerning EnCluDL, we notice that the δ -FWER is not
 644 controlled when $n = 50$ except for $C = 100$. This is not surprising since the δ -FWER
 645 control is asymptotic and $n = 50$ is not sufficient. In terms of support recovery, the
 646 problem gets easier with larger n , but only EnCluDL benefits strongly from an increase
 647 of n . Third, we investigate the influence of the level of correlation between neighboring
 648 covariates (ρ). Regarding FWER control, desparsified lasso does not control the FWER
 649 when $\rho = 0.5$. Regarding the statistical power of EnCluDL, as one would expect, when
 650 the spatial structure is strong *i.e.*, $\rho > 0.9$, it is relevant to pick larger clusters, *i.e.*, to
 651 take a smaller C . Indeed, to make a relevant choice for C , data structure has to be taken
 652 into account to derive good covariates' clustering; this is true up to the limit given by
 653 the desired spatial tolerance. A good clustering creates clusters that are weakly corre-
 654 lated and contains covariates that are highly correlated. This observation is linked to
 655 assumption (ii) of Prop. 4.1 or to assumption (ii.a) and (ii.b) of Prop. B.1. Finally,
 656 we consider the results for different support sizes coded by the active region width h .
 657 Sparsity is a crucial assumption for desparsified lasso and then for EnCluDL. Also, when
 658 p (or C) increases the required sparsity is greater. This explains why when $h = 8$ and
 659 $C \geq 300$, the empirical δ -FWER is slightly above the expected nominal rate. Regarding
 660 the statistical power of EnCluDL, as one could expect, when the active regions are large,
 661 it is relevant to use large clusters. However, it can be difficult to estimate this parameter
 662 in advance, thus we prefer to consider desired spatial tolerance parameter δ and data
 663 structure to set C .

664 E Proofs

665 E.1 Proof of Prop. 4.1 and Prop. B.1

666 First, we start by the proof of Prop. 4.1 which is derived from Bühlmann et al. [2013,
 667 Proposition 4.3]:

Proof. With assumption (ii) and Bühlmann et al. [2013, Proposition 4.3], we have, for
 all $c \in [C]$:

$$\boldsymbol{\theta}_c^* = |G_c| \sum_{j \in G_c} w_j \boldsymbol{\beta}_j^* ,$$

where, for all $j \in G_c$:

$$w_j = \frac{\sum_{k \in G_c} \Sigma_{j,k}}{\sum_{k \in G_c} \sum_{k' \in G_c} \Sigma_{k,k'}} .$$

668 From assumption (i), we have $w_j > 0$ for all $j \in G_c$. Assumption (iii) ensures that, for
 669 all $j \in G_c$, the β_j^* have the same sign. Then, θ_c^* is of the same sign as the β_j^* and is
 670 non-zero only if there exists $j \in G_c$ such that $\beta_j^* \neq 0$. \square

671 Now, we give the proof of **Prop. B.1** which is mainly derived from **Bühlmann et al.**
 672 **[2013, Proposition 4.4]**:

Proof. With assumption (ii.a) and (ii.b) and **Bühlmann et al. [2013, Proposition 4.4]**, we have, for all $c \in [C]$:

$$\theta_c^* = |G_c| \sum_{j \in G_c} w'_j \beta_j^* + \kappa_c ,$$

where

$$w'_j = \frac{\sum_{k \in G_c} \text{Cov}(\mathbf{X}_{\cdot,j}, \mathbf{X}_{\cdot,k} \mid \{\mathbf{Z}_{\cdot,c'} : c' \neq c\})}{\sum_{k \in G_c} \sum_{k' \in G_c} \text{Cov}(\mathbf{X}_{\cdot,k}, \mathbf{X}_{\cdot,k'} \mid \{\mathbf{Z}_{\cdot,c'} : c' \neq c\})} ,$$

and, for all $c \in [C]$

$$|\kappa_c| \leq (\nu_c / \tau_c) \|\beta^*\|_1 .$$

Let us define $\tilde{\theta}$ by

$$\tilde{\theta}_c = |G_c| \sum_{j \in G_c} w'_j \beta_j^* .$$

Then,

$$\theta^* = \tilde{\theta} + \kappa ,$$

673 And, similarly as in the proof of **Prop. 4.1**, from assumption (i) and (iii), $\tilde{\theta}_c$ is of the
 674 same sign as the β_j^* for $j \in G_c$ and is non-zero only if there exists $j \in G_c$ such that
 675 $\beta_j^* \neq 0$. \square

676 E.2 Proof of **Prop. 4.2**

Before going through the proof of **Prop. 4.2**, we introduce the grouping function g that matches the covariate index to its corresponding group index:

$$\begin{aligned} g : [p] &\rightarrow [C] \\ j &\mapsto c \quad \text{if } j \in G_c . \end{aligned}$$

Then, (7) can be rewritten as follows:

$$\begin{aligned} \text{for all } j \in [p], \quad \hat{p}_j &= \hat{p}_{g(j)}^{\mathcal{G}} , \\ \text{for all } j \in [p], \quad \hat{q}_j &= \hat{q}_{g(j)}^{\mathcal{G}} . \end{aligned} \tag{18}$$

Proof. (i) Suppose that we are under $H_0^\delta(j)$. Since the cluster diameters are all smaller than δ , all the covariates in $G_{g(j)}$ have a corresponding weight equal to zero. Thus, using [Prop. 4.1](#), we have $\boldsymbol{\theta}_{g(j)}^* = 0$, *i.e.*, we are under $H_0(G_{g(j)})$. Under this last null-hypothesis, using [\(11\)](#) and [\(18\)](#), we have:

$$\text{for all } \alpha \in (0, 1), \mathbb{P}(\hat{p}_{g(j)}^{\mathcal{G}} \leq \alpha) = \mathbb{P}(\hat{p}_j \leq \alpha) = \alpha .$$

677 This last result being true for any $j \in N^\delta$, we have shown that the elements of the family
678 \hat{p} control the δ -type 1 error.

(ii) As mentioned in [Sec. 4.4](#), we know that, the family $\hat{q}^{\mathcal{G}}$ controls the FWER, *i.e.*, for $\alpha \in (0, 1)$ we have $\mathbb{P}(\min_{c \in N_{\mathcal{G}}} \hat{q}_c^{\mathcal{G}} \leq \alpha) \leq \alpha$. Let us denote by $g^{-1}(N_{\mathcal{G}})$ the set of indexes of covariates that belong to the groups of $N_{\mathcal{G}}$, *i.e.*, $g^{-1}(N_{\mathcal{G}}) = \{j \in [p] : g(j) \in N_{\mathcal{G}}\}$. Again, given that all the cluster diameters are smaller than δ and using [Prop. 4.1](#), if $j \in N^\delta$ then $g(j) \in N_{\mathcal{G}}$. That is to say $N^\delta \subset g^{-1}(N_{\mathcal{G}})$. Then, we have:

$$\min_{j \in N^\delta} (\hat{q}_j) \geq \min_{j \in g^{-1}(N_{\mathcal{G}})} (\hat{q}_j) .$$

We can also notice that:

$$\begin{aligned} \min_{j \in g^{-1}(N_{\mathcal{G}})} (\hat{q}_j) &= \min_{j \in g^{-1}(N_{\mathcal{G}})} (\hat{q}_{g(j)}^{\mathcal{G}}) \\ &= \min_{g(j) \in N_{\mathcal{G}}} (\hat{q}_{g(j)}^{\mathcal{G}}) . \end{aligned}$$

Replacing $g(j) \in [C]$ by $c \in [C]$, and using [\(6\)](#), we obtain:

$$\text{for all } \alpha \in (0, 1), \mathbb{P}(\min_{j \in N^\delta} (\hat{q}_j) \leq \alpha) \leq \mathbb{P}(\min_{c \in N_{\mathcal{G}}} \hat{q}_c^{\mathcal{G}} \leq \alpha) \leq \alpha .$$

679 This last result states that the family $(\hat{q}_j)_{j \in [p]}$ controls the δ -FWER. □

680 **E.3 Proof of [Prop. 4.3](#)**

681 The proof of [Prop. 4.3](#) is inspired by the one proposed by [Meinshausen et al. \[2009\]](#).
682 However, it is subtly different since we can not remove the term $\min_{j \in N^\delta}$ and have to
683 work with it to obtained the desired inequality. First, we start by making a short remark
684 about the γ -quantile quantity.

Definition E.1 (empirical γ -quantile). *For a set V of real numbers and $\gamma \in (0, 1)$, let*

$$\gamma\text{-quantile}(V) = \min \left\{ v \in V : \frac{1}{|V|} \sum_{w \in V} \mathbf{1}_{w \leq v} \geq \gamma \right\} . \quad (19)$$

Remark E.1. *For a set of real number V and for a $a \in \mathbb{R}$, let us define the quantity $\pi(a, V)$ by the following:*

$$\pi(a, V) = \frac{1}{|V|} \sum_{v \in V} \mathbf{1}(v \leq a) \quad (20)$$

685 Then, for $\gamma \in (0, 1)$, the two events $E_1 = \{\pi(a, V) \geq \gamma\}$ and $E_2 = \{\gamma\text{-quantile}(V) \leq a\}$
686 are identical.

Now, we give the proof of [Prop. 4.3](#).

Proof. First, one can notice that, from [\(8\)](#), we have:

$$\min_{j \in N^\delta}(\tilde{q}_j(\gamma)) \geq \min \left\{ 1, \gamma\text{-quantile} \left(\left\{ \min_{j \in N^\delta} \left(\frac{\hat{q}_j^{(b)}}{\gamma} \right) : b \in [B] \right\} \right) \right\} .$$

Then, for $\alpha \in (0, 1)$:

$$\begin{aligned} \mathbb{P} \left(\min_{j \in N^\delta}(\tilde{q}_j(\gamma)) \leq \alpha \right) &\leq \mathbb{P} \left(\min \left\{ 1, \gamma\text{-quantile} \left(\left\{ \min_{j \in N^\delta} \left(\frac{\hat{q}_j^{(b)}}{\gamma} \right) : b \in [B] \right\} \right) \right\} \leq \alpha \right) \\ &= \mathbb{P} \left(\gamma\text{-quantile} \left(\left\{ \min_{j \in N^\delta} \left(\frac{\hat{q}_j^{(b)}}{\gamma} \right) : b \in [B] \right\} \right) \leq \alpha \right) . \end{aligned}$$

Using [Rem. E.1](#), for $\gamma \in (0, 1)$, with:

$$V = \left\{ \min_{j \in N^\delta} \left(\frac{\hat{q}_j^{(b)}}{\gamma} \right) : b \in [B] \right\} \quad \text{and} \quad a = \alpha ,$$

and noticing that:

$$\pi \left(\alpha, \left\{ \min_{j \in N^\delta} \left(\frac{\hat{q}_j^{(b)}}{\gamma} \right) : b \in [B] \right\} \right) = \frac{1}{B} \sum_{b=1}^B \mathbb{1} \left\{ \min_{j \in N^\delta}(\hat{q}_j^{(b)}) \leq \alpha \gamma \right\} ,$$

then, we have:

$$\mathbb{P} \left(\gamma\text{-quantile} \left(\left\{ \min_{j \in N^\delta} \left(\frac{\hat{q}_j^{(b)}}{\gamma} \right) : b \in [B] \right\} \right) \leq \alpha \right) = \mathbb{P} \left(\frac{1}{B} \sum_{b=1}^B \mathbb{1} \left\{ \min_{j \in N^\delta}(\hat{q}_j^{(b)}) \leq \alpha \gamma \right\} \geq \gamma \right) .$$

Then, the Markov inequality gives:

$$\mathbb{P} \left(\frac{1}{B} \sum_{b=1}^B \mathbb{1} \left\{ \min_{j \in N^\delta}(\hat{q}_j^{(b)}) \leq \alpha \gamma \right\} \geq \gamma \right) \leq \frac{1}{\gamma} \mathbb{E} \left[\frac{1}{B} \sum_{b=1}^B \mathbb{1} \left\{ \min_{j \in N^\delta}(\hat{q}_j^{(b)}) \leq \alpha \gamma \right\} \right] .$$

Then, using the assumption that the B families $(\hat{q}_j^{(b)})_{j \in [p]}$ control of the δ -FWER (last inequality), we have:

$$\frac{1}{\gamma} \mathbb{E} \left[\frac{1}{B} \sum_{b=1}^B \mathbb{1} \left\{ \min_{j \in N^\delta}(\hat{q}_j^{(b)}) \leq \alpha \gamma \right\} \right] = \frac{1}{\gamma} \frac{1}{B} \sum_{b=1}^B \mathbb{P} \left(\min_{j \in N^\delta}(\hat{q}_j^{(b)}) \leq \alpha \gamma \right) \leq \alpha .$$

Finally, we have shown that, for $\alpha \in (0, 1)$:

$$\mathbb{P} \left(\min_{j \in N^\delta}(\tilde{q}_j(\gamma)) \leq \alpha \right) \leq \alpha .$$

688 This establishes that the family $(\tilde{q}_j(\gamma))_{j \in [p]}$ controls the δ -FWER. \square

689 **E.4 Proof of Theorem 4.1**

690 To show Theorem 4.1, we connect the previous results: Prop. 4.1, Prop. 4.2 and Prop. 4.3.
691 First, we prove that clustered inference algorithms produce a p-value family that controls
692 the δ -FWER.

693 *Proof.* Assuming the noise model (1), assuming that Ass. 2.1 and Ass. 2.2 are verified
694 for a distance parameter larger than δ and that the clustering diameter is smaller than
695 δ , then we directly obtain the assumption (i) and (iii) of Prop. 4.1. This means that the
696 compressed representation has the correct pattern of non-zero coefficients, in particular
697 it does not include in the support clusters of null-only covariates. Additionally, if one is
698 able to perform a valid statistical inference on the compressed model (3), *i.e.*, to produce
699 cluster-wise p-values such that (4) holds, then Prop. 4.2 ensures that the p-value family
700 constructed using the de-grouping method presented in (7) controls the δ -FWER. \square

701 Now, we prove that ensembled clustered inference algorithms produce a p-value fam-
702 ily that controls the δ -FWER.

703 *Proof.* Given the above arguments, the p-value families produced by clustered inference
704 algorithms subject to all clusterings fulfilling the theorem hypotheses control the δ -
705 FWER. Then, using the aggregation method given by (8), we know from Prop. 4.3 that
706 the aggregated p-value family also controls the δ -FWER. \square