

## $\ell_0$ -PENALIZED MAXIMUM LIKELIHOOD FOR SPARSE DIRECTED ACYCLIC GRAPHS

BY SARA VAN DE GEER AND PETER BÜHLMANN

*ETH Zürich*

We consider the problem of regularized maximum likelihood estimation for the structure and parameters of a high-dimensional, sparse directed acyclic graphical (DAG) model with Gaussian distribution, or equivalently, of a Gaussian structural equation model. We show that the  $\ell_0$ -penalized maximum likelihood estimator of a DAG has about the same number of edges as the minimal-edge I-MAP (a DAG with minimal number of edges representing the distribution), and that it converges in Frobenius norm. We allow the number of nodes  $p$  to be much larger than sample size  $n$  but assume a sparsity condition and that any representation of the true DAG has at least a fixed proportion of its non-zero edge weights above the noise level. Our results do not rely on the faithfulness assumption nor on the restrictive strong faithfulness condition which are required for methods based on conditional independence testing such as the PC-algorithm.

**1. Introduction.** Directed acyclic graphs (DAGs) and corresponding directed graphical models are key concepts for causal inference, see for example the books by Spirtes, Glymour and Scheines (2000) and Pearl (2000). From an estimation point of view, a first step consists in estimating the Markov equivalence class of the true underlying causal DAG based on observational data, and from this, one can infer identifiable causal effects and lower bounds for all causal effects (Maathuis, Kalisch and Bühlmann, 2009). This strategy has been applied to, and to a certain extent validated using high-throughput, and hence high-dimensional, data in biology (Maathuis et al., 2010). It is of primary importance to understand limitations and potential of methods in terms of subtle and often uncheckable assumptions, and in this respect, our results here shed new light.

We focus here on the problem of estimating the Markov equivalence class of DAGs, or more generally of a so-called minimal-edge I-MAP, in the setting of observational Gaussian data where the number  $p$  of variables or

---

*AMS 2000 subject classifications:* Primary 62F12; secondary 62F30

*Keywords and phrases:* Causal inference, Faithfulness condition, Gaussian structural equation model, Graphical modeling, High-dimensional inference

nodes in the DAG may greatly exceed sample size  $n$ . We consider the  $\ell_0$ -penalized maximum likelihood estimator, and we relate and compare our new results and conditions to the popular PC-algorithm (Spirtes, Glymour and Scheines, 2000) and its theoretical analysis. To the best of our knowledge, the latter is so far the only work providing theoretical guarantees for inferring the Markov equivalence class of DAGs in the high-dimensional setting. We emphasize that the popular  $\ell_1$ -norm regularization of the likelihood is inappropriate here, leading to an objective function which is not constant over equivalent DAGs encoding the same distribution. On the other hand,  $\ell_0$ -penalization leads to invariant scores over equivalent DAGs. The computational difficulties are primarily due to the non-convex constraint that the directed graph should be acyclic, and the additional issue with the  $\ell_0$ -in comparison to e.g.  $\ell_1$ -norm penalization is, surprisingly, rather in favor of the former. A computationally feasible algorithm for exact  $\ell_0$ -penalized maximum likelihood estimation for the Markov equivalence class of DAGs has been proposed by Silander and Myllymäki (2006); for larger graphs, with more than about 20 nodes, approximate algorithms can be used (Chickering, 2002; Hauser and Bühlmann, 2012). More details are given in Section 2.1.

1.1. *Relation to other work.* We analyze the  $\ell_0$ -penalized maximum likelihood estimator for the equivalence class of DAGs in the Gaussian setting. Pioneering work for the low-dimensional case on this problem has been done by Chickering (2002) who proved consistency of the BIC-score and provided an algorithm, called greedy equivalent search (GES), which greedily proceeds in the space of Markov equivalence classes. While the GES-algorithm can also be used in the high-dimensional scenario (Hauser and Bühlmann, 2012), the asymptotic consistency of BIC is established only for the case with a fixed distribution (with  $p < \infty$ ) where the sample size  $n \rightarrow \infty$ . Chickering's first significant analysis does not provide any insights for the high-dimensional case with its subtle interplay of signal strength, noise level and identifiability conditions.

Robins et al. (2003) present refined analyses for causal inference under the view point of uniform consistency as sample size  $n \rightarrow \infty$ . There, problematic issues with the so-called faithfulness condition (see Section 1.2) arise, and Zhang and Spirtes (2003) introduce the notion of strong faithfulness (see (2)), as a way to address some of the raised major problems. None of these works consider high-dimensional inference, but their pointing to the faithfulness condition and its version are important.

Kalisch and Bühlmann (2007) provide consistency results of the PC-algorithm (Spirtes, Glymour and Scheines, 2000) for estimating the Markov equivalence class of DAGs based on Gaussian observational data, in the high-dimensional, sparse setting. One of the conditions used is a restricted version of strong faithfulness in (2). Our analysis with the penalized MLE is completely different and circumvents faithfulness and strong faithfulness conditions which are often very restrictive as shown by Uhler et al. (2012), see also below.

The theoretical high-dimensional analysis presented here is very different and more challenging than for multivariate regression or covariance estimation, due to the *unknown* order among the variables. For known order, as e.g. in time series problem, Shojaie and Michailidis (2010) present results for estimation of high-dimensional DAGs; but the case with unknown order considered in the present paper requires major new theoretical ideas and development.

1.2. *Directed graphical and structural equation models.* Consider the following model. There is a DAG  $D_0$  whose  $p$  nodes correspond to random variables  $X_1, \dots, X_p$ : assume that

$$(1) \quad \begin{aligned} X_1, \dots, X_p &\sim \mathcal{N}_p(0, \Sigma_0) \text{ with density } f_{\Sigma_0}(\cdot), \\ \mathcal{N}_p(0, \Sigma_0) &\text{ is Markovian with respect to } D_0, \end{aligned}$$

where the Markov property can be understood as the factorization property where the joint Gaussian density  $f_{\Sigma_0}(x_1, \dots, x_p)$  can be factorized as

$$f_{\Sigma_0}(x_1, \dots, x_p) = \prod_{j=1}^p f_{\Sigma_0}(x_j | x_{\text{pa}(j)}),$$

with  $\text{pa}(j)$  denoting the set of parents of node  $j$  (Lauritzen, 1996, cf.).

It is well-known that in general, there exists another DAG  $D$  such that the distribution  $\mathcal{N}(0, \Sigma_0)$  is Markovian with respect to  $D$ . Assuming faithfulness (see below), the set of all such other DAGs build the so-called Markov equivalence class  $\mathcal{E}(D_0)$  which can be characterized in terms of a bi-directed graph (Andersson, Madigan and Perlman, 1997, cf.). The Markov equivalence class  $\mathcal{E}(D_0)$  can be identified from the observational data distribution  $\mathcal{N}(0, \Sigma_0)$  under the assumption of faithfulness.

*Definition of faithfulness (Spirtes, Glymour and Scheines, 2000, cf.):* For a DAG  $D$ , a distribution  $P$  is called *faithful* with respect to  $D$  if and only if all conditional independences are encoded by the DAG  $D$ .

Faithfulness is stronger than a Markov assumption: the latter allows to infer some conditional independences from the DAG while the former requires that *all* of them can be inferred from the DAG (i.e. only the ones which are entailed by a Markov condition). Failure of faithfulness is “rare”, having Lebesgue measure zero, if the edge weights (the coefficients in the equivalent linear structural equation model) are chosen from a distribution which is absolutely continuous with respect to Lebesgue measure. However, for statistical estimation, we often require sufficiently strong detectability of conditional dependencies, given by the notion of strong faithfulness.

*Definition of strong faithfulness in the Gaussian case (Zhang and Spirtes, 2003):* For a DAG  $D$ , a Gaussian distribution  $P$  is called  $\tau$ -strongly faithful with respect to  $D$  if and only if

$$(2) \quad \min\{|\text{Parcorr}(X_j, X_k|X_S)|; \text{Parcorr}(X_j, X_k|X_S) \neq 0, \\ S \subseteq \{1, \dots, p\} \setminus \{j, k\}, j, k \in \{1, \dots, p\} (j \neq k)\} \geq \tau.$$

A typical requirement is strong faithfulness with  $\tau \asymp \sqrt{\text{sparsity} \cdot \log(p)/n}$ , see also below for the PC-algorithm. Strong faithfulness can be viewed as a condition of “signal strength” in terms of non-zero partial correlations. As shown in Uhler et al. (2012), strong faithfulness is a very restrictive condition for many DAGs, and the same applies to a slightly weaker restricted strong faithfulness assumption (Uhler et al., 2012, cf.), see also Section 4.3.2. At the same time, it is essentially unavoidable for any algorithm for inferring the Markov equivalence class  $\mathcal{E}(D_0)$  which relies on conditional independence testing. The most prominent example is the PC-algorithm (Spirtes, Glymour and Scheines, 2000): consistent estimation for the Markov equivalence class  $\mathcal{E}(D_0)$  is proved in Kalisch and Bühlmann (2007) for the Gaussian case assuming a strong faithfulness condition. The results in this paper for the  $\ell_0$ -penalized MLE do not require a strong faithfulness condition as in (2) (nor the slightly weaker restricted strong faithfulness condition): the reason is that the method is *not* relying on conditional independence testing but rather on penalized parameter estimation in terms of a linear structural equation model as explained next.

A Gaussian DAG model as in (1) can always be equivalently represented as a linear structural equation model:

$$(3) \quad X_j = \sum_{k \in \text{pa}(j)} \beta_{kj}^0 X_k + \epsilon_j \quad (j = 1, \dots, p),$$

where  $\epsilon_1, \dots, \epsilon_n$  are independent,  $\epsilon_j \sim \mathcal{N}(0, |\omega_j^0|^2)$  and  $\epsilon_j$  independent of  $\{X_k; k \in \text{pa}(j)\}$ ; note that  $\text{pa}(j) = \text{pa}_{D_0}(j)$  depends on the true DAG  $D_0$ .

**2. The setting and the estimator.** We use here and in the sequel a terminology which does not rely on the standard language from graphical modeling since the required basics for the Gaussian case (see models (1) and (3)) can be developed in a straightforward way.

We consider  $n$  i.i.d. observations from the structural equation model (3) which is equivalent to model (1). We denote by  $X := (X_1, \dots, X_p)$  the  $n \times p$  data matrix with  $n$  i.i.d. rows, each of them being  $\mathcal{N}(0, \Sigma_0)$ -distributed, where  $\Sigma_0$  is a non-singular covariance matrix. The relations between the variables in a row can be represented as

$$(4) \quad X = XB_0 + E,$$

where  $B_0 := (\beta_{k,j}^0)$  is a  $p \times p$  matrix with  $\beta_{j,j}^0 = 0$  for all  $j$ , and where  $E$  as an  $n \times p$  matrix of noise vectors  $E := (\epsilon_1, \dots, \epsilon_p)$ , with  $\epsilon_j$  independent of  $X_k$  whenever  $\beta_{k,j}^0 \neq 0$ .<sup>1</sup> Furthermore,  $E$  has  $n$  i.i.d. rows which are  $\mathcal{N}(0, \Omega_0)$ -distributed, with  $\Omega_0 := \text{diag}(|\omega_1^0|^2, \dots, |\omega_p^0|^2)$  a  $p \times p$  diagonal matrix.

The model (4) implies that

$$\Sigma_0 = [(I - B_0)^{-1}]^T \Omega_0 [(I - B_0)^{-1}].$$

We call  $(B_0, \Omega_0)$  a DAG corresponding to  $\Sigma_0$ .<sup>2</sup> The set of edges of this DAG is denoted by  $s_0 := s_{B_0} := \{(k, j) : \beta_{k,j}^0 \neq 0\}$ , and in fact,  $\beta_{k,j}^0 \neq 0$  encodes for a directed edge  $k \rightarrow j$ . We will assume in Condition 3.4 that  $(B_0, \Omega_0)$  is sparse, in the sense that its number of (directed) edges  $s_0$  is small.

As described in Section 1.2 with the concept of a Markov equivalence class, there are several DAGs  $(\tilde{B}_0, \tilde{\Omega}_0)$  describing the same  $\Sigma_0$  and thus the same Gaussian distribution  $P = \mathcal{N}(0, \Sigma_0)$ . Throughout this paper,  $(B_0, \Omega_0)$  is defined as a DAG with a minimal number of edges (and it may not be unique). We call such a DAG a minimal-edge I-MAP.<sup>3</sup>

<sup>1</sup>Note that in relation to the true DAG  $D_0$  in model (3),  $\beta_{k,j}^0 = 0$  for  $k \notin \text{pa}_{D_0}(j)$ . We do not make such explicit constraints here since we aim for a smallest DAG representing the distribution of  $X$ .

<sup>2</sup>This deviates from the classical definition where the DAG is only a (directed acyclic) graph; we use the short terminology ‘‘DAG’’ for the whole graphical model with the distribution and the graph encoded by the coefficient matrix  $B$  and the error variances  $\Omega$ .

<sup>3</sup>It is a minimal I-MAP (Spirtes, Glymour and Scheines, 2000, Sec.2.3.1.) with the additional property that it has minimal number of edges.

**Remark 2.1** We call two DAGs  $(B_1, \Omega_1)$  and  $(B_2, \Omega_2)$  equivalent if  $\Theta(B_1, \Omega_1) = \Theta(B_2, \Omega_2)$  and if in addition they have the same number of edges.<sup>4</sup> In our analysis, we will identify DAGs which are in the same equivalence class. Thus, our aim is to estimate an arbitrary member of the equivalence class of  $(B_0, \Omega_0)$ , by a suitable member in the equivalence class of  $(\hat{B}, \hat{\Omega}_0)$ .

2.1. *The  $\ell_0$ -penalized maximum likelihood estimator.* We use a penalized maximum likelihood procedure to estimate the DAG  $(B_0, \Omega_0)$ . Let

$$\Sigma_n := X^T X / n$$

be the empirical covariance matrix based on the observations  $X$ . Given a  $p \times p$  non-singular covariance matrix  $\Sigma$ , with inverse  $\Theta := \Sigma^{-1}$ , the minus log-likelihood is proportional to

$$l_n(\Theta) := \text{trace}(\Theta \Sigma_n) - \log \det(\Theta).$$

We consider inverse covariance matrices that can be represented as a DAG. That is, we let

$$\Theta := \Theta(B, \Omega) := (I - B)\Omega^{-1}(I - B)^T,$$

where  $(B, \Omega)$  is a DAG. The latter means that  $\Omega$  is a positive diagonal matrix and that, up to a permutation  $\pi$  of the rows and columns,  $B$  can be written as a lower-diagonal matrix.

The  $\ell_0$ -penalized maximum likelihood estimator is

$$(5) \quad \hat{\Theta} = \arg \min_{B, \Omega} \left\{ l_n(\Theta) + \lambda^2 s_B : \Theta = \Theta(B, \Omega), \text{ for some DAG } (B, \Omega) \text{ with } B \in \mathcal{B} \right\}.$$

Here  $s_B$  is the number of non-zero elements in  $B$  (corresponding to the number of edges in the DAG) and  $\lambda \geq 0$  is a tuning parameter. The estimator is denoted by  $\hat{\Theta} := \Theta(\hat{B}, \hat{\Omega})$ . It has  $\hat{s} := s_{\hat{B}}$  edges. The collection  $\mathcal{B}$  is the set of all edge weights  $B$  of DAGs  $(B, \Omega)$  which have at most  $\alpha n / \log p$  incoming edges (parents) at each node, where  $\alpha > 0$  is given (see Condition 3.3), or a subset thereof. We will discuss this restriction in Subsection 4.2.

---

<sup>4</sup>This definition is not the same as for a Markov equivalence class. Assuming faithfulness, both definitions coincide.

We throughout assume  $B_0 \in \mathcal{B}$ , i.e. that the restrictions one puts on the edge weights are correct.

The  $\ell_0$ -penalty in the estimator ensures that the penalized likelihood remains the same among all equivalent representations, e.g., among all DAGs from the same Markov equivalence class or among the equivalence class described in Remark 2.1 above. This would not be true when choosing for example an  $\ell_1$ -norm penalization  $\|B\|_1 := \sum_{k,j} |\beta_{k,j}|$ . From a computational point of view, the main difficulty is the optimization over the space of DAGs  $\mathcal{B}$ : current algorithms are tailored for the  $\ell_0$ -penalty (see next paragraph), and in this sense, optimization of the  $\ell_0$ -penalized log-likelihood is better tractable than using an  $\ell_1$ -norm regularization. For problems with up to about  $p \approx 20$  nodes, dynamic programming can be used (Silander and Myllymäki, 2006), while for large-scale applications, greedy equivalent search has been reported to perform well (Chickering, 2002; Hauser and Bühlmann, 2012).

The dynamic programming method (Silander and Myllymäki, 2006) is optimizing the  $\ell_0$ -penalized negative log-likelihood in (5) over the space of all DAGs with  $B \in \mathcal{B}$ . It crucially relies on the assumption that the objective function, called the score, can be decomposed locally such that  $\text{score}(D) = \sum_{j=1}^p g(X_j, X_{\text{pa}_D(j)})$  for some function  $g(\cdot)$ , where  $g(X_j, X_{\text{pa}_D(j)})$  involves only data from the variables  $j$  and  $\text{pa}_D(j)$ . The  $\ell_0$ -penalized score is decomposable, whereas say  $\ell_1$ -norm penalization does not lead to a decomposable score. The greedy equivalent search algorithms (Chickering, 2002; Hauser and Bühlmann, 2012) are greedily optimizing the objective function in (5) in the space of equivalence classes, a much smaller space than the space of DAGs. The greedy steps are forward, backward and turning edges moves such that the score is improved most when stepping from one Markov equivalence class to another one: this can be done very efficiently without enumerating all DAG members in the equivalence class. Such greedy equivalent search algorithms rely crucially on the fact that the objective score function is constant within an equivalence class and that the score is decomposable: again, this is true with  $\ell_0$ -penalization but fails for  $\ell_1$ -norm regularization.

The previous paragraph already gives a good reason why  $\ell_0$ -penalization is to be preferred over the  $\ell_1$ -norm analogue, namely, that the computational techniques are tailored and surprisingly easier for the former than the latter ( $\ell_1$ -norm regularization would require to search over the whole space of DAGs, and the dynamic programming trick mentioned above cannot be used). Furthermore,  $\ell_1$ -norm regularization is usually motivated as a convex relaxation and this is not true in our context since the DAG constraint in (5) for the matrix  $B \in \mathcal{B}$  remains to cause that the optimization problem is

non-convex. In addition, the price to pay, in a standard regression problem with an  $\ell_1$ -norm regularization, is a bias which should be further controlled with e.g. adaptive  $\ell_1$ -norm regularization (Zou, 2006) or thresholding (van de Geer, Bühlmann and Zhou, 2011). We have not considered a theoretical analysis of the  $\ell_1$ -norm penalized maximum likelihood estimator of a DAG. We believe this not to be more difficult than the  $\ell_0$ -norm penalized maximum likelihood estimator although the proofs (see Subsection 7.1 for an outline) may need different arguments.

2.1.1. *The main results and their implications.* We show in Theorem 3.1 that with a choice of the tuning parameter  $\lambda^2$  of order  $\log p/n((p/s_0)\vee 1)$ , the number of edges  $\hat{s}$  of the estimator is of the same order of magnitude as the number of edges  $s_0$  of a true underlying minimal-edge I-MAP. Moreover, we show that  $\hat{B}$  and  $\hat{\Omega}$  converges in Frobenius norm to some  $\tilde{B}_0$  and  $\tilde{\Omega}_0$  respectively, where  $\Theta(\tilde{B}_0, \tilde{\Omega}_0) = \Theta_0$  is a representation of the true DAG with  $\tilde{s}$  edges (see Section 2.2), and with  $\tilde{s}$  of the same order of magnitude as  $s_0$ . The rate of convergence is of order  $\lambda^2 s_0$ . To arrive at this result, we need that at least a fixed proportion of the non-zero coefficients of any representation of the true DAG is above the “noise level”, the latter being of order  $\sqrt{\log p/n}(\sqrt{p/s_0}\vee 1)$  (see Condition 3.5): in analogy to regression, we call this the “beta-min” condition. Some of our other conditions are trivially satisfied when  $p = \mathcal{O}(n/\log(n))$  is sufficiently small.

The “noise level” indicates two regimes for  $(n, p, s_0)$ . If  $s_0$  is at least of the order of  $p$  (or larger), than the “noise level” is of the order  $\sqrt{\log(p)/n}$  which is small even if  $p$  is very large relative to  $n$ . This scenario is often realistic saying that a fixed non-zero proportion of the nodes has at least one parent: we call it the standard sparsity regime. The other scenario is for  $s_0 \ll p$ , corresponding to very sparse DAGs, which we call the ultra-high sparsity regime. The reason for the two different noise level scenarios is that we estimate  $p$  error variances in  $\Omega_0$ : when  $s_0 \ll p$ , the term from estimating these error variances is dominating.

When making a more stringent “beta-min” condition and choosing the regularization parameter of larger order than the “noise level” (or the same order if  $s_0 = \mathcal{O}(1)$ ), we can recover the minimal-edge I-MAP.

The  $\ell_0$ -penalized MLE can be easily adapted to the case where the noise variances in  $\Omega_0$  are known up to a scalar, for example when all noise variances are known to be equal but their value is unknown. Then, the true DAG can be identified from the distribution (Peters and Bühlmann, 2012). We show that in this case, under an identifiability condition on the noise variances,

the  $\ell_0$ -penalized maximum likelihood estimator finds the representation (and hence the true DAG) with the prescribed noise variances, and the rate of convergence for the Frobenius norm is of order  $\lambda^2 s_0$  (see Theorem 5.1). We assume in this context that  $p$  is sufficiently smaller than  $n/\log n$ .

**Remark 2.2** *If the minimal-edge I-MAP can be inferred and assuming in addition the faithfulness condition (but not requiring strong faithfulness), see Section 1.2, we can recover the true underlying Markov equivalence class. This then allows to derive bounds for causal inference, exactly along the lines of Maathuis, Kalisch and Bühlmann (2009).*

2.2. *Permutations and the order of the variables.* The model (4) can be written as

$$X_j = X\beta_j^0 + \epsilon_j, \quad j = 1, \dots, p,$$

with  $\beta_j^0$  the  $j$ -th column of  $B_0$ . Let us write for any vector  $\beta \in \mathbb{R}^p$ ,

$$\|X\beta\|^2 := \beta^T \Sigma_0 \beta, \quad \|X\beta\|_n^2 := \beta^T \Sigma_n \beta.$$

For a permutation  $\pi$  of  $\{1, \dots, p\}$ , which plays the role of an order of the variables, we let  $\tilde{B}_0(\pi)$  be the matrix obtained by doing a Gram-Schmidt orthogonalization for  $\|\cdot\|$ , starting with  $X_{\pi_p}$  and finishing by projecting  $X_{\pi_1}$  on  $X_{\pi_2}, \dots, X_{\pi_p}$ . Moreover, we let  $\tilde{\Omega}_0(\pi)$  be the diagonal matrix of the error variances. Note thus that  $\tilde{B}_0(\pi)$  is lower-diagonal after permutation of its rows and columns. Furthermore,

$$\Theta_0 = \Theta_0(\tilde{B}_0(\pi), \tilde{\Omega}_0(\pi)), \forall \pi.$$

The set of incoming edges at node  $j$  (non-zero coefficients of the  $j$ th column of  $\tilde{B}_0(\pi)$ ) is denoted by  $\tilde{S}_j(\pi)$ , and we let  $\tilde{s}_j(\pi) := |\tilde{S}_j(\pi)|$ . Moreover, we let  $\tilde{s}(\pi) = \sum_{j=1}^p \tilde{s}_j(\pi)$  be the total number of edges of  $\tilde{B}_0(\pi)$ . Thus,  $\tilde{s}(\pi) = \tilde{s}_{\tilde{B}_0(\pi)}$ .

Let  $\hat{B}$  - or one of the members in its equivalence class - be lower-diagonal after permutation  $\hat{\pi}$ , and define  $\tilde{B}_0 := B_0(\hat{\pi})$ . The number of edges of  $\tilde{B}_0$  is denoted by  $\tilde{s} = \tilde{s}(\hat{\pi})$ . The DAGs  $(\hat{B}, \hat{\Omega})$  and  $(\tilde{B}_0, \tilde{\Omega}_0)$  share the same lower-diagonal structure (but not necessarily the same set of zero coefficients). We will show that  $\tilde{s} := \tilde{s}(\hat{\pi})$  is with large probability of the same order of magnitude as  $s_0$  (see Theorem 3.1). Thus, if the true DAG  $(B_0, \Omega_0)$  is sparse, then with large probability the DAG  $(\tilde{B}_0(\hat{\pi}), \tilde{\Omega}_0(\hat{\pi}))$  is sparse as well, which means that on average, the number of incoming edges at a node is small.

Note that  $\hat{\pi}$  is a random permutation and that there are in total a large number (namely  $p!$ ) permutations. Analytical control over these  $p!$  permutations requires a very different technique than dealing with known order (Shojaie and Michailidis, 2010) or with multivariate regression or covariance estimation problems. We explain this in more detail in Section 7.4.1.

2.2.1. *AR(1)-model as an example.* Suppose the true DAG is a directed chain from  $X_p$  along  $X_{p-1}, \dots, X_2$  to  $X_1$  with a corresponding structural equation model (AR(1)-model):

$$\begin{aligned} X_p &= \epsilon_p, \\ X_j &= \beta^0 X_{j+1} + \epsilon_j \quad (j = 1, \dots, p-1), \end{aligned}$$

where  $\epsilon \sim \mathcal{N}(0, \Omega_0)$  with  $\Omega_0 = \text{diag}(1 - (\beta^0)^2, \dots, 1 - (\beta^0)^2, 1)$  and  $|\beta^0| < 1$ . The error variances are chosen such that  $\text{Var}(X_j) = 1$  for all  $j$ . The covariance matrix is of Toeplitz form  $(\Sigma_0)_{ij} = (\beta^0)^{|i-j|}$  and the model satisfies the directed global Markov property which is equivalent to the concept of d-separation (Lauritzen, 1996, cf. Sec.3.2.2).

Therefore, we have that projecting

$$X_{\pi_j} \text{ on } X_{\pi_{j+1}}, \dots, X_{\pi_p} \quad (j = 1, \dots, p-1)$$

leads to at most two non-zero regression coefficients in every column of  $\tilde{B}_0(\pi)$  (corresponding to the largest index  $k_1 < \pi_j$  and smallest index  $k_2 > \pi_j$  if  $\pi_{j+1}, \dots, \pi_p$  contains indices smaller and larger than  $\pi_j$ ; or corresponding to the largest  $k < \pi_j$  if  $\pi_{j+1}, \dots, \pi_p$  contains only smaller indices than  $\pi_j$ ; or corresponding to the smallest  $k > \pi_j$  if  $\pi_{j+1}, \dots, \pi_p$  contains only larger indices than  $\pi_j$ ). Thus, we have that  $\tilde{s}_j(\pi) \leq 2$  for all  $j$  and all  $\pi$  and hence Condition 3.4, given below, holds.

The absolute values of the non-zero coefficients in  $\tilde{B}_0(\pi) = (\tilde{\beta}_{k,j}^0(\pi))$  decrease monotonely as the index-distance  $d(j) = \min_{k=j+1, \dots, p} |\pi_k - \pi_j|$  increases. Thus, for fixed  $j$  and whenever  $d(j) > \Delta$  for some (large) value of  $\Delta$ , there are at most two (since  $\tilde{s}_j(\pi) \leq 2$ ) coefficients with  $|\tilde{\beta}_{k,j}^0(\pi)| \leq C(\Delta)$  for some value  $C(\Delta)$  (which decreases as  $\Delta$  increases and which depends on  $\beta^0$ ). Therefore, clearly, there are at most  $2(\lfloor p/\Delta \rfloor + 1)$  coefficients (edges) whose values satisfy  $|\tilde{\beta}_{k,j}^0(\pi)| \leq C(\Delta)$ , and all other non-zero coefficients (at least  $p - \lfloor p/\Delta \rfloor - 2$ )<sup>5</sup> are larger than  $C(\Delta)$ . For e.g.  $\Delta = 3$  this implies that

<sup>5</sup>There are at least  $p - (\lfloor p/\Delta \rfloor + 1)$  indices (nodes)  $j$  with  $d(j) \leq \Delta$ ; and there is at least one non-zero coefficient (edge) from all of them except one (the starting node). The value  $C = C(\Delta)$  can be chosen appropriately, for any fixed  $\Delta$ .

there are at most  $2(\lfloor p/3 \rfloor + 1)$  edges with non-zero coefficients being smaller than  $C(\Delta = 3)$ , and at least  $p - \lfloor p/3 \rfloor - 2$  edges with non-zero coefficients larger than  $C(\Delta = 3)$ . This implies that Condition 3.5, given below, holds with a value  $C(\Delta = 3)$  of order 1.

**3. Conditions and main result.** We write  $\Sigma_0 =: (\sigma_{k,j})$  and we let  $\sigma_j^2 := \sigma_{j,j}$ ,  $j = 1, \dots, p$ .

**Condition 3.1** For some constant  $\sigma_0^2$ , it holds that

$$\max_{1 \leq j \leq p} \sigma_j^2 \leq \sigma_0^2.$$

**Condition 3.2** The smallest eigenvalue  $\Lambda_{\min}^2$  of  $\Sigma_0$  is non-zero. (See also (6)).

**Condition 3.3** For a given constant  $\alpha > 0$ , it holds that for any  $B = (\beta_1, \dots, \beta_p) \in \mathcal{B}$ , where  $\mathcal{B}$  is as in (5), that  $s_{\beta_j} \leq \alpha n / \log p$  for all  $j = 1, \dots, p$ , where for a vector  $\beta \in \mathbb{R}^p$  we denote the cardinality of its support set by  $s_\beta := \#\{\beta_k \neq 0\}$ .

**Condition 3.4** For some constant  $\tilde{\alpha}$  and any permutation  $\pi$ , and all  $j$ ,

$$\tilde{s}_j(\pi) \leq \tilde{\alpha} n / \log p,$$

where  $\tilde{s}_j(\pi) = s_{\tilde{\beta}_j^0}$  is the number of incoming edges of the DAG  $(\tilde{B}_0(\pi), \tilde{\Omega}_0(\pi))$  at node  $j$ . (See also (7).)

**Condition 3.5** There exist constants  $0 \leq \eta_1 < 1$  and  $0 < \eta_0^2 < 1 - \eta_1$ , such that for any permutation  $\pi$ , the DAG  $(\tilde{B}_0(\pi), \tilde{\Omega}_0(\pi))$  (which has  $\tilde{s}(\pi)$  edges) has at least  $(1 - \eta_1)\tilde{s}(\pi)$  edges  $(k, j)$  with  $|\tilde{\beta}_{k,j}^0(\pi)| > \sqrt{\log p/n}(\sqrt{p/s_0} \vee 1)/\eta_0$ .

Following Bühlmann and van de Geer (2011), we refer to Condition 3.5 as the “beta-min” condition. It is a “kind of replacement” of the strong faithfulness condition in (2) that is required for consistency of the PC algorithm and variants thereof, see Section 1.2. A detailed discussion about the assumptions is given in Section 4.

In the current section, we will present an asymptotic formulation for clarity. We will provide a non-asymptotic result in Section 7. Our results depend on  $\Sigma_0$  via the constants  $\sigma_0$ ,  $\Lambda_{\min}$ , on the further constants  $\gamma_0 := (\alpha, \tilde{\alpha}, \eta_0, \eta_1)$  used in the conditions, and on  $\alpha_0$  where  $1 - \alpha_0$  is the confidence level of the statement.

We assume that we can take  $\gamma_0$  sufficiently small. Moreover, we state the results with  $\alpha_0 := (4/p) \wedge .05$  to avoid digressions concerning the confidence level. Explicit expressions can be found in Theorem 7.4, where we simplified the situation by assuming  $n$  is sufficiently large and  $\log p/n$  is sufficiently small.

With the notation  $z = \mathcal{O}(1)$  we mean that  $z$  can be bounded by a constant depending only on  $\sigma_0$  and  $\Lambda_{\min}$ . Moreover, with  $z \asymp 1$  we mean  $z = \mathcal{O}(1)$  and  $1/z = \mathcal{O}(1)$ . Furthermore, the Frobenius norm is defined as  $\|B\|_F = (\sum_{j,k=1}^p |\beta_{k,j}|^2)^{1/2}$  for a  $p \times p$  matrix  $B$  with elements  $\beta_{k,j}$ .

**Theorem 3.1** *Assume Conditions 3.1, 3.2, 3.3, 3.4 and 3.5, with  $\gamma_0 := (\alpha, \tilde{\alpha}, \eta_0, \eta_1)$  sufficiently small, but allowing  $1/\|\gamma_0\|_1 = \mathcal{O}(1)$ . Let  $1 - \alpha_0$  be the confidence level, with  $\alpha_0 := (4/p) \wedge .05$ . Then for a choice*

$$\lambda^2 \asymp \frac{\log p}{n} \left( \frac{p}{s_0} \vee 1 \right),$$

*it holds that with probability at least  $1 - \alpha_0$ ,*

$$\|\hat{B} - \tilde{B}_0\|_F^2 + \|\hat{\Omega} - \tilde{\Omega}_0\|_F^2 = \mathcal{O}(\lambda^2 s_0),$$

*where  $(\tilde{B}_0, \tilde{\Omega}_0)$  are defined in Section 2.2, and*

$$\hat{s} \asymp \tilde{s} \asymp s_0.$$

The proof is given in Section 7. Theorem 7.4 gives some explicit bounds.

**Remark 3.1** *Note that if the true permutation  $\pi_0$  defined by  $\tilde{B}_0(\pi_0) = B_0$  is known, then from (multiple) regression theory the optimal rate of convergence in Frobenius norm will be of order  $(p + s_0) \log p/n$  (with  $p \log p/n$  being a lower bound due to estimating the  $p$  residual variances). Hence, Theorem 3.1 says that this rate can also be achieved when not knowing  $\pi_0$ . As in a multiple regression setup, a natural normalization of the Frobenius norm is to divide it by  $p$ . With this normalized norm, the estimator is consistent when the average number of incoming edges  $s_0/p$  is of small order  $n/\log p$ .*

**Remark 3.2** *If the beta-min condition (Condition 3.5) holds with  $\eta_1 = 0$  and with very small values for  $\eta_0 := \eta_0(\pi)$ , namely of order  $1/\tilde{s}(\pi)$ , then one obtains the screening property: all edges in  $(\tilde{B}_0, \tilde{\Omega}_0)$  are then with large probability also present in  $(\hat{B}, \hat{\Omega})$ .*

*Moreover, by taking  $\lambda^2 := \lambda^2(s_0)$  very large (of order  $s_0 \log p/n$ ), one can obtain with large probability that  $\hat{s} \leq s_0$ . In other words, by imposing a*

strong beta-min condition, which is severe if  $s_0$  is large, one recovers with high probability the edges of the minimal-edge I-MAP exactly. However, in Theorem 3.1, we do not use such values for  $\eta_0$  and  $\lambda$ , but instead  $\lambda^2 \asymp \log p/n$  (when  $p = \mathcal{O}(s_0)$ ) and  $\eta_0 \asymp 1$ . Thus, we generally do not recover the true edges. This is the price for dealing with a large  $p$  situation and an  $s_0$  possibly growing in  $n$ . Such problems do not show up in asymptotics with  $p$  fixed.

#### 4. A discussion of the conditions.

4.1. *Bounds for the noise variances.* For all  $\pi$  and  $j$  and for any  $\beta_j$  with  $\beta_{j,j} = 0$ , we have  $\|X_j - X\beta_j\| = \|X\beta_j^-\|$  where  $\beta_{k,j}^- = -\beta_{k,j}$  for  $k \neq j$  and  $\beta_{j,j}^- = 1$ . It follows that for any  $\pi$  and  $j$ ,

$$|\tilde{\omega}_j^0(\pi)|^2 = \|X_j - X\tilde{\beta}_j^0(\pi)\|^2 \geq \Lambda_{\min}^2,$$

with  $\Lambda_{\min}^2$  the smallest eigenvalue of  $\Sigma_0$ . Moreover, clearly  $|\tilde{\omega}_j^0(\pi)|^2 \leq \sigma_j^2$ . Hence, Conditions 3.1 and 3.2 imply that for all  $\pi$  and  $j$

$$0 < \Lambda_{\min}^2 \leq |\tilde{\omega}_j^0(\pi)|^2 \leq \sigma_0^2.$$

Furthermore,  $\Lambda_{\min}^2 > 0$  is implied by

$$(6) \quad \min_j |\omega_j^0|^2 > 0,$$

since  $\det(\Sigma_0) = \det(\Omega_0) = \prod_{j=1}^p \omega_j^0$ . Thus, Condition 3.2 is equivalent to (6).

4.2. *Overfitting.* Condition 3.3 will ensure that the penalized minus log-likelihood cannot become minus infinity. If  $n$  or more edges are allowed at a node, say at node  $j$ , the estimator will overfit the data at this node, giving a residual variance  $\hat{\omega}_j^2 = 0$ . The penalized minus log-likelihood is proportional to  $\sum_{j=1}^p \log \hat{\omega}_j^2 + \lambda^2 \hat{s}$  which will be  $-\infty$  if one allows that  $\hat{\omega}_j$  vanishes. Note that the penalty as such does not prevent this type of overfitting. Therefore, we need a restriction on the class of possible DAGs, and Condition 3.3 serves this purpose. We will show in Lemma 7.5 that Conditions 3.1, 3.2 and 3.3 imply that for an appropriate constant  $K_0 > 0$ , it holds for all  $j$  that  $\hat{\omega}_j \geq 1/K_0$  with large probability.

4.3. *The beta-min condition.* One may circumvent the beta-min condition if one allows for edges with weights below some noise level  $\lambda_*$  to be set to zero. Here,  $\lambda_* := \sqrt{\log p/n}/\eta_0^*$  for some suitable  $\eta_0^* > 0$ . Instead of trying to estimate the true DAG  $(B_0, \Omega_0)$ , one now aims at estimating its best sparse approximation  $(B_0^*, \Omega_0^*)$ , which is defined as follows. Let for any DAG  $(B, \Omega)$ , and for  $\Theta = \Theta(B, \Omega)$ , the weights  $B_\Theta(\pi)$  be obtained by doing the Gram-Schmidt orthogonalization for  $\|\cdot\|_\Sigma$ , where  $\Sigma = \Theta^{-1}$ , and  $\|X\beta\|_\Sigma^2 := \beta^T \Sigma \beta$ ,  $\beta \in \mathbb{R}^p$ . Thus  $B_\Theta(\pi)$  is lower-diagonal after the permutation  $\pi$  of its rows and columns, and for appropriate  $\Omega_\Theta(\pi)$ , the DAG  $(B_\Theta(\pi), \Omega_\Theta(\pi))$  satisfies

$$\Theta = \Theta(B_\Theta(\pi), \Omega_\Theta(\pi)).$$

Let  $s_\Theta(\pi) := s_{B_\Theta(\pi)}$  be the number of edges of  $B_\Theta(\pi)$ . Connecting this with our previous notation, we note that

$$B_{\Theta_0}(\pi) = \tilde{B}_0(\pi), \quad \Omega_{\Theta_0}(\pi) = \tilde{\Omega}_0(\pi), \quad s_{\Theta_0}(\pi) = \tilde{s}(\pi).$$

Let now for some constant  $\eta_0^* > 0$ ,

$$s_\Theta^*(\pi) := \#\left\{ (k, j) : |\beta_{\Theta, k, j}(\pi)| > \sqrt{\log p/n}(\sqrt{p/s_\Theta(\pi)} \vee 1)/\eta_0^* \right\}.$$

We then take

$$\Theta_0^* := \arg \min \{ l(\Theta) : \Theta = (B, \Omega) \text{ a DAG, } s_\Theta^*(\pi) \geq (1 - \eta_1^*)s_\Theta(\pi), \forall \pi \},$$

where  $0 \leq \eta_1^* < 1$  and  $l(\Theta) = \text{trace}(\Theta \Sigma_0) - \log \det(\Theta) = \mathbb{E} l_n(\Theta)$  is the theoretical counterpart of the minus log-likelihood. (Note that  $\Theta_0 := \Sigma_0^{-1}$  is the overall minimizer of  $l(\Theta)$ .) We let  $(B_0^*, \Omega_0^*)$  be a solution of

$$\Theta_0^* = \Theta_0^*(B_0^*, \Omega_0^*)$$

with the minimum number of edges. With constants  $\eta_0^*$  and  $\eta_1^*$  sufficient small, one may replace  $\Theta_0 = \Theta_0(B_0, \Omega_0)$  by  $\Theta_0^* = \Theta_0^*(B_0^*, \Omega_0^*)$  in our analysis. In this way, one can avoid the beta-min condition, provided that the bias term that will now appear in the bounds is small enough.

4.3.1. *The beta-min condition and the number of edges.* We further note that Conditions 3.1, 3.2 and 3.5 imply Condition 3.4 with

$$(7) \quad \tilde{\alpha} = \frac{\sigma_0^2 \eta_0^2}{\Lambda_{\min}^2 (1 - \eta_1)}.$$

This is because for all  $j$ ,

$$\frac{(1 - \eta_1) \tilde{s}_j(\pi) \log p}{\eta_0^2 n} \leq \|\tilde{\beta}_j^0(\pi)\|_2^2 \leq \|X \tilde{\beta}_j^0(\pi)\|^2 / \Lambda_{\min}^2 \leq \sigma_0^2 / \Lambda_{\min}^2.$$

4.3.2. *The strong beta-min condition in comparison to strong faithfulness.* The beta-min condition as stated in Condition 3.5 is rather weak. In order to make a comparison with strong faithfulness, which focuses on exact edge recovery, we consider the stronger version as discussed in Remark 3.2. As written in this remark, recovery of a minimal-edge I-MAP is guaranteed with a value for the lower bound on the weights of the non-zero edges of the order  $s_0 \sqrt{\log(p)/n} = p \sqrt{\log(p)/n}$  assuming  $s_0 \asymp p$ : such a value in the beta-min condition is reasonable in the regime  $p = o(\sqrt{n/\log(n)})$ .

Although this seems rather restrictive at first sight, it is often better than the situation for the PC-algorithm which requires restricted strong faithfulness (Uhler et al., 2012, cf.) for consistent estimation of the Markov equivalence class. There, the best dimensionality range is achieved for bounded-degree trees which restricts  $p = o(\sqrt{n/\log(n)})$  (Uhler et al., 2012, Sec.5.1) to the same order of magnitude as above while for other graphs the constraint on  $p$  can be much stronger (see below).

The AR(1) model in Section 2.2.1 is an example where the beta-min condition holds with a value of order 1. Regarding restricted strong faithfulness, the model is a bounded degree tree leading to a dimensionality restriction  $p = o(\sqrt{n/\log(n)})$  which is the same as for our (rough) result about edge recovery with the  $\ell_0$ -penalized MLE. We can extend the analysis to an AR( $k$ ) model with fixed  $k \geq 2$ . Then, using analogous reasoning as for the AR(1) in Section 2.2.1, the beta-min condition holds for a value of order 1 still leading to a dimensionality range  $p = o(\sqrt{n/\log(n)})$ . For restricted faithfulness to hold, however, a more severe constraint for consistent estimation of the Markov equivalence class comes into play since the corresponding graph is not a tree anymore. We use the technique in Uhler et al. (2012) and evaluate another polynomial than for the considered cases there: for example, for an AR(2) model, the dimensionality restriction becomes  $o((n/\log(n))^{1/10})$ .

We will discuss in Section 5 the case where the error variances are the same, i.e.,  $\Omega_0 = \omega^0 I$ . We then only need a beta-min condition for the true underlying DAG instead of all permutations (see the discussion after Theorem 5.1). Thus, for the scenario  $p = o(\sqrt{n/\log(n)})$  in the equal variance case, the beta-min condition is very reasonable for any DAG. This is in sharp contrast to the constraint arising from restricted strong faithfulness: if the underlying DAG is say a certain bipartite graph, the corresponding dimensionality for consistent edge recovery becomes  $p = o(\log(n))$  (Uhler et al., 2012, Sec.5.1).

Finally, we note that when focusing on bounding false positive selections as in Theorem 3.1, the  $\ell_0$ -penalized MLE is justified for the  $p \gg n$  setting.

4.4. *The high sparsity regime where  $s_0 \ll p$ .* The reason why we see a term  $p/s_0 \vee 1$  appearing in the tuning parameter  $\lambda^2$  (see Theorem 3.1) and in the beta-min condition (Condition 3.5) is due to the estimation of the  $p$  unknown variances, which gives a term of order  $p \log p/n$  in our bounds for the squared Frobenius norm. If  $s_0 \ll p$ , the true DAG has many disconnected components, and in fact it then has many isolated points. The variables in one component are uncorrelated with those in another component. We see this in the zeroes in the matrix  $\Sigma_0$ . The connected components and isolated points are easily detected by  $\Sigma_n$ , assuming that non-zero correlations are at least  $\sqrt{\log p/n}/\eta_c$  in absolute value for an appropriate (sufficiently small) constant  $\eta_c$ . Then we can do the analysis connected component by connected component. To summarize, the situation  $p = \mathcal{O}(s_0)$  appears to be the most interesting. Alternatively, when the noise variances  $\Omega_0$  are known up to a scalar (for example if it is known that all noise variances are equal), we need not estimate these variances anymore and the term of order  $p \log p/n$  does not appear in the bounds, provided an identifiability condition on the noise variances holds and  $p$  is sufficiently smaller than  $n/\log n$ . This will be shown in the next section.

**5. The case of equal variances.** Suppose that the noise variances  $\{|\omega_j^0|^2\}_{j=1}^p$  are known up to a multiplicative scalar. To simplify the exposition, let us assume that

$$\omega_1^0 = \dots = \omega_p^0 = 1.$$

The  $\ell_0$ -penalized maximum likelihood estimator now becomes

$$(8) \quad \hat{B} := \arg \min \left\{ \text{trace} \left( (I - B)(I - B)^T \Sigma_n \right) + \lambda^2 s_B : \right. \\ \left. (B, I) \text{ a DAG, } B \in \mathcal{B} \right\},$$

where  $B$  is as in (5).

The main Theorem 3.1 as well as the remarks in Section 3.2 apply to the estimator (8) as well, assuming exactly the same Conditions 3.1 - 3.5.

For the case where  $p = \mathcal{O}(n/\log(n))$  is sufficiently small, we obtain consistent estimation of the true underlying DAG and we gain in comparison to the main Theorem 3.1 by excluding the additional factor  $(p/s_0 \vee 1)$ . We make the following assumptions.

**Condition 5.1** *There exists a constant  $\eta_\omega > 0$  such that for all  $\tilde{\Omega}_0(\pi) \neq I$ ,*

$$\frac{1}{p} \sum_{j=1}^p \left( |\tilde{\omega}_j^0(\pi)|^2 - 1 \right)^2 > 1/\eta_\omega.$$

**Condition 5.2** *There exists a constant  $\alpha_*$  such that*

$$p \leq \alpha_* n / \log n.$$

We call Condition 5.1 the “omega-min” condition. It leads to identification of the DAG with equal variances. Condition 5.2 ensures that the rate of convergence is fast enough to ensure that eventually we choose the right permutation. Note that it implies Conditions 3.3 and 3.4, with  $\alpha = \tilde{\alpha} = \alpha_*$ .

Let  $\pi_0$  be defined by  $\tilde{B}_0(\pi_0) = B_0$ . Since  $B_0$  is identifiable from the observational distribution  $\mathcal{N}(0, \Sigma_0)$  (Peters and Bühlmann, 2012), see also Section 2.1.1,  $\pi_0$  corresponds to the unique true ordering of the variables.

**Theorem 5.1** *Assume Conditions 3.1 and 3.2, and Conditions 5.1 and 5.2. Let  $\alpha_0 := (4/p) \wedge .05$ . Then for  $\gamma_* := (\alpha_*, \eta_\omega)$  suitably small, but allowing  $1/\|\gamma_*\|_1 = \mathcal{O}(1)$ , and for  $\lambda^2 \asymp \log p/n$ , it holds with probability at least  $1 - \alpha_0$ , that  $\hat{\pi} = \pi_0$ , and*

$$\|\hat{B} - B_0\|_F^2 + \lambda^2 \hat{s} = \mathcal{O}(\lambda^2 s_0).$$

The proof is given in Subsection 7.6.

Thus, we find  $\hat{s} = \mathcal{O}(s_0)$ , but we do not show  $\hat{s} \asymp s_0$ . To establish the latter, one again needs a beta-min condition, but this time only on the DAG  $(B_0, I)$ , and not on any of the other representations  $(\tilde{B}_0(\pi), \tilde{\Omega}_0(\pi))$  with  $\pi \neq \pi_0$ . This is a much simplified and weaker assumption than in Condition 3.5. Furthermore, choosing  $\lambda^2 \asymp s_0 \log p/n$  sufficiently large, exact edge recovery follows by the beta-min condition for the true DAG  $(B_0, I)$ , that is, the condition that  $\min\{|\beta_{k,j}^0| : \beta_{k,j} \neq 0\} > s_0 \sqrt{\log p/n}/\eta^0$  for some sufficiently small  $\eta^0 > 0$ .

5.1. *The non-Gaussian case.* To avoid technical digressions in our proofs, we assume a Gaussian distribution for the observations where zero correlations mean independence. We use in Lemma 7.4 that if for some  $\tilde{\epsilon}_j$ ,  $\mathbb{E}\tilde{\epsilon}_j = 0$ , then also the conditional expectation of  $\tilde{\epsilon}_j$  given variables  $X_k$  that are uncorrelated with  $\tilde{\epsilon}_j$  is zero. In the non-Gaussian case, this is no longer true. However, one can still derive similar results, along a line of proof that does

not use conditioning but instead concentration inequalities for averages of products of random variables (empirical covariances). This means that our results go through for observations which are sub-Gaussian. The proofs then rely on concentration inequalities of Bernstein-type. The main adjustments of our proofs are then as follows. We assume that the rows of  $X$  form an i.i.d. sequence of sub-Gaussian vectors as defined in Loh and Wainwright (2012) and replace Theorem 7.3 by their Lemma 15. In Lemma 7.4 we assume  $\epsilon$  and  $Z$  are sub-Gaussian and uncorrelated, and replace the empirical squared norm  $\|X\beta\|_n^2 := [\sum_{i=1}^n (X\beta)_i^2/n]$  by the theoretical squared norm  $\|X\beta\|^2$ . We can then apply similar arguments as used in Lemma 15 of Loh and Wainwright (2012). In Lemma 7.1, we no longer use the empirical squared norm but instead the theoretical one. Theorem 7.2 needs virtually no adjustments.

**6. Conclusions.** We establish the first results of the  $\ell_0$ -penalized MLE for estimation of the minimal-edge I-MAP (the smallest DAG which can generate the data-generating distribution) in the high-dimensional sparse setting. Thereby, we avoid the faithfulness condition, and the strong faithfulness assumption (2) or its restricted version (Uhler et al., 2012, cf.); the latter is necessary for consistency of the PC-algorithm (Spirtes, Glymour and Scheines, 2000). The (restricted) strong faithfulness condition is typically very strong (Uhler et al., 2012) and hence, our results contribute in relaxing such very restrictive assumptions.

Our main assumption is the beta-min Condition 3.5 (which implies the sparsity Condition 3.4, see Section 4.3.1): as an example, the AR(1)-model in Section 2.2.1 fulfills it, even if  $p \gg n$ . The noise level is of the order  $\sqrt{\log(p)/n(p/s_0 \vee 1)}$ : the additional factor  $(p/s_0 \vee 1)$  occurs due to estimation of  $p$  variances in  $\Omega_0$ . However, the interesting scenario is for the case where  $s_0 \geq \text{const.}p$  since  $s_0 \ll p$  corresponds to a DAG where most nodes are isolated having no edges to other nodes; thus, for  $s_0 \geq \text{const.}p$ , we obtain the usual noise level of the order  $\sqrt{\log(p)/n}$ , as in high-dimensional regression problems.

For the equal variance case with  $p = \mathcal{O}(n/\log(n))$  sufficiently small, our result in Theorem 5.1 (and its comment below) is most clear in that we essentially only require the beta-min Condition 3.5 for the true DAG  $B_0$ , i.e., a substantially relaxed assumption, and the identifiability Condition 5.1 for the error variances: we can then recover the true underlying unique DAG  $B_0$ . Thus, we have identified an important class of models where estimation of the order of variables and the true underlying DAG is possible without requiring the badly limiting (restricted) strong faithfulness condition (2).

## 7. Proofs.

7.1. *A brief outline of the proofs.* We first consider the proof of Theorem 3.1 which treats the case of unknown variances  $\{(\omega_j^0)^2\}$ . In Lemma 7.1 of Subsection 7.2, we present a bound for  $\sum_{j=1}^p [(\hat{\omega}_j^0)^2/\hat{\omega}_j^2 - 1]^2$  and for  $\sum_{j=1}^p \|X(\hat{\beta}_j - \tilde{\beta}_j^0)\|_n^2$  using the empirical norm  $\|v\|_n := [\sum_{i=1}^n v_i^2/n]^{1/2}$ ,  $v \in \mathbb{R}^n$ . The result follows from a straightforward manipulation of likelihoods, but it is assumed there that one is on the part of the probability space where the random components behave well. The study of these random components is deferred to Subsections 7.4.1, 7.4.2 and 7.4.3. First, the bound of Lemma 7.1 is refined because it involves the number of edges  $\tilde{s}$  of the DAG formed by using the random permutation  $\hat{\pi}$  that the penalized maximum likelihood estimator chooses. Subsection 7.3 presents the tools to deal with this by exploiting the beta-min condition. The idea here is that if the Frobenius norm between  $\hat{B}$  and  $\tilde{B}$  is small, the number of edges of  $\tilde{B}$  cannot be much larger than those of  $\hat{B}$ .

A substantial part of the proof of Theorem 3.1 goes into showing that with large probability we are on a set of the form  $\cap_{k=0}^3 \mathcal{T}_k$  where the random components behave well. Let us first discuss  $\mathcal{T}_1$ . Here, a uniform inequality holds for the empirical correlation between the projections and error terms in a Gram Schmidt orthogonalization. For a fixed permutation  $\pi$  it is rather standard to control these empirical correlations. The new element is that we have to control them uniformly over all permutations  $\pi$  in order to show that  $\mathcal{T}_1$  has large probability. We do this in Subsection 7.4.1, where the arguments used are explained just before Theorem 7.1. In the set  $\mathcal{T}_2$  all empirical variances of the error terms in a Gram Schmidt orthogonalization are close to their expectations. We show in Subsection 7.4.2 that uniformly over all  $\pi$  this is true with large probability. The set  $\mathcal{T}_3$  gives bounds for  $\|\beta\|_2$  in terms of  $\|X\beta\|_n$  and the number of non-zero coefficients in  $\beta$ . We show in Theorem 7.3 that  $\mathcal{T}_3$  has large probability. This makes it possible to move from empirical norms to Frobenius norms and moreover shows that with large probability the  $\{\hat{\omega}_j^2\}$  are bounded away from zero. The latter event is defined as the set  $\mathcal{T}_0$ .

For the proof of Theorem 5.1 where the variances  $(\omega_j^0)^2$  are all known to be equal to one, we use the same structure. We assume that we are on the set  $\cap_{k=1}^3 \mathcal{T}_k$ , and use straightforward manipulations of likelihoods.

7.2. *Bounds on a subset of the probability space.* We present some explicit bounds assuming we are on a set of the form  $\cap_{k=0}^3 \mathcal{T}_k$ , where the sets  $\mathcal{T}_k$

are defined below. Then we show in Subsections 7.4.1, 7.4.2 and 7.4.3 that each  $\mathcal{T}_k$ ,  $k = 0, \dots, 3$  has large probability for an appropriate choice of the constants and of the parameters  $\lambda_1$ ,  $\lambda_2$  and  $\lambda_3$  involved in the definition of these sets. In fact, we will show that one can take  $\lambda_1 \asymp \lambda_2 \asymp \lambda_3 \asymp \sqrt{\log p/n}$ .

Let for some constant  $K_0 > 0$ ,

$$\mathcal{T}_0 := \{\hat{\omega}_j^2 \geq 1/K_0^2, \forall j\}.$$

Let us write  $X_k \perp \tilde{\epsilon}_j$  if  $X_k$  and  $\tilde{\epsilon}_j$  are independent. For all  $\pi$  and  $j$ , define  $\tilde{\epsilon}_j(\pi) = X_j - X\tilde{\beta}_j^0(\pi)$ , and  $\tilde{\mathcal{B}}_j(\pi) := \{\beta_j : X_k \perp \tilde{\epsilon}_j(\pi), \forall \beta_{k,j} \neq 0\}$ . Moreover, let  $\tilde{\mathcal{B}}(\pi) := \{B = (\beta_1, \dots, \beta_p) \in \mathcal{B} : \beta_j \in \tilde{\mathcal{B}}_j(\pi) \forall j\}$ . For some  $\delta_1 > 0$  and some  $\lambda_1 > 0$ , write

$$\begin{aligned} \mathcal{T}_1 := & \left\{ 2 \sum_{j=1}^p |\tilde{\epsilon}_j^T(\pi) X(\beta_j - \tilde{\beta}_j^0(\pi))|/n \leq \delta_1 \sum_{j=1}^p \|X(\beta_j - \tilde{\beta}_j^0(\pi))\|_n^2 \right. \\ & \left. + \lambda_1^2(s + \tilde{s}(\pi))/\delta_1, \forall B = (\beta_1, \dots, \beta_p) \in \tilde{\mathcal{B}}(\pi), \forall \pi \right\}, \end{aligned}$$

We let for some  $\lambda_2 > 0$ ,

$$\mathcal{T}_2 := \left\{ \sum_{j=1}^p \left( \frac{\|\tilde{\epsilon}_j(\pi)\|_n^2 - |\tilde{\omega}_j^0(\pi)|^2}{|\tilde{\omega}_j^0(\pi)|^2} \right)^2 \leq 4\lambda_2^2(p + \tilde{s}(\pi)), \forall \pi \right\},$$

where we recall the notation  $\|v\|_n^2 := v^T v/n$ ,  $v \in \mathbb{R}^n$ . Finally, for some some  $\delta_3 > 0$  and some  $\lambda_3 > 0$ , let  $\mathcal{T}_3$  be the set

$$\mathcal{T}_3 := \left\{ \|X\beta\|_n \geq \left[ \delta_3 - \lambda_3 \sqrt{s_\beta} \right] \|\beta\|_2, \forall \beta \right\}.$$

Recall that  $s_\beta := \#\{\beta_k \neq 0\}$ .

**Lemma 7.1** *Define  $(\tilde{B}_0, \tilde{\Omega}_0) := (B_0(\hat{\pi}), \Omega_0(\hat{\pi}))$  and  $\tilde{s} := s_{\tilde{B}_0}$ . Assume that Condition 3.1 holds. Suppose we are on  $\cap_{k=0}^2 \mathcal{T}_k$  with  $0 < \delta_1 < 1/K_0^2$  and  $0 < \delta_2 < 1/(2K_0^4\sigma_0^4)$ . Take the tuning parameter  $\lambda^2 > \lambda_1^2/\delta_1 + \lambda_2^2/\delta_2$ . Then*

$$\begin{aligned} & \left( \frac{1}{K_0^2} - \delta_1 \right) \sum_{j=1}^p \|X(\hat{\beta}_j - \tilde{\beta}_j^0)\|_n^2 + \left( \frac{1}{2K_0^4\sigma_0^4} - \delta_2 \right) \sum_{j=1}^p \left( \frac{\hat{\omega}_j^2 - |\tilde{\omega}_j^0|^2}{|\hat{\omega}_j|^2} \right)^2 \\ & + \left( \lambda^2 - \frac{\lambda_1^2}{\delta_1} - \frac{\lambda_2^2}{\delta_2} \right) \hat{s} \leq \lambda^2 s_0 + \frac{\lambda_2^2(p + \tilde{s})}{\delta_2} + \frac{\lambda_1^2 \tilde{s}}{\delta_1}. \end{aligned}$$

**Proof.** Let  $\tilde{\epsilon} := \tilde{\epsilon}(\hat{\pi})$ . We apply the Basic Inequality

$$l_n(\hat{\Theta}) + \lambda^2 \hat{s} \leq l_n(\Theta_0) + \lambda^2 s_0,$$

or equivalently

$$p + \sum_{j=1}^p \log \hat{\omega}_j^2 + \lambda^2 \hat{s} \leq \sum_{j=1}^p \frac{\|\epsilon_j\|_n^2}{|\omega_j^0|^2} + \sum_{j=1}^p \log |\omega_j^0|^2 + \lambda^2 s_0.$$

which gives, using  $\log(\det(\Sigma_0)) = \sum_{j=1}^p \log |\omega_j^0|^2 = \sum_{j=1}^p \log |\tilde{\omega}_j^0|^2$ ,

$$\sum_{j=1}^p \log \left( \frac{\hat{\omega}_j^2}{|\tilde{\omega}_j^0|^2} \right) + \lambda^2 \hat{s} \leq \sum_{j=1}^p \left( \frac{\|\epsilon_j\|_n^2}{|\omega_j^0|^2} - 1 \right) + \lambda^2 s_0.$$

Since  $\hat{\omega}_j^2 \geq 1/K_0^2$  (since we are on  $\mathcal{T}_0$ ) and  $|\tilde{\omega}_j^0|^2 \leq \sigma_0^2$  (by Condition 3.1), we know that  $|\tilde{\omega}_j^0|^2/\hat{\omega}_j^2 \leq K_0^2 \sigma_0^2$ . But then, using  $\log(1+x) \leq x - x^2/(2(1+c)^2)$ ,  $-1 < x \leq c$ , we get

$$\log \left( \frac{\hat{\omega}_j^2}{|\tilde{\omega}_j^0|^2} \right) = -\log \left( \frac{|\tilde{\omega}_j^0|^2}{|\hat{\omega}_j^2} \right) \geq -\left( \frac{|\tilde{\omega}_j^0|^2}{|\hat{\omega}_j^2} - 1 \right) + \frac{1}{2K_0^4 \sigma_0^4} \left( \frac{|\tilde{\omega}_j^0|^2}{|\hat{\omega}_j^2} - 1 \right)^2.$$

We plug this back into the Basic Inequality to get

$$\sum_{j=1}^p \frac{\hat{\omega}_j^2 - |\tilde{\omega}_j^0|^2}{|\hat{\omega}_j^2} + \frac{1}{2K_0^4 \sigma_0^4} \left( \frac{\hat{\omega}_j^2 - |\tilde{\omega}_j^0|^2}{|\hat{\omega}_j^2} \right)^2 + \lambda^2 \hat{s} \leq \sum_{j=1}^p \left( \frac{\|\epsilon_j\|_n^2}{|\omega_j^0|^2} - 1 \right) + \lambda^2 s_0.$$

Rewrite this to

$$\begin{aligned} & \sum_{j=1}^p \frac{\|X(\hat{\beta}_j - \tilde{\beta}_j^0)\|_n^2}{\hat{\omega}_j^2} + \frac{1}{2K_0^4 \sigma_0^4} \sum_{j=1}^p \left( \frac{\hat{\omega}_j^2 - |\tilde{\omega}_j^0|^2}{|\hat{\omega}_j^2} \right)^2 + \lambda^2 \hat{s} \\ & \leq 2 \sum_{j=1}^p \frac{\tilde{\epsilon}_j^T X(\hat{\beta}_j - \tilde{\beta}_j^0)/n}{\hat{\omega}_j^2} + \sum_{j=1}^p \left( \frac{\|\epsilon_j\|_n^2}{|\omega_j^0|^2} - 1 \right) - \sum_{j=1}^p \left( \frac{\|\tilde{\epsilon}_j\|_n^2 - |\tilde{\omega}_j^0|^2}{|\hat{\omega}_j^2} \right) + \lambda^2 s_0. \end{aligned}$$

We now apply

$$\begin{aligned} & \sum_{j=1}^p \left( \frac{\|\tilde{\epsilon}_j\|_n^2 - |\tilde{\omega}_j^0|^2}{|\hat{\omega}_j^2} \right) \\ & = \sum_{j=1}^p \left( \frac{\|\tilde{\epsilon}_j\|_n^2 - |\tilde{\omega}_j^0|^2}{|\tilde{\omega}_j^0|^2} \right) + \sum_{j=1}^p \left( \frac{\|\tilde{\epsilon}_j\|_n^2 - |\tilde{\omega}_j^0|^2}{|\tilde{\omega}_j^0|^2} \right) \left( \frac{|\tilde{\omega}_j^0|^2 - \hat{\omega}_j^2}{\hat{\omega}_j^2} \right). \end{aligned}$$

But, by the Cauchy-Schwarz inequality and using that we are on  $\mathcal{T}_2$ ,

$$\begin{aligned} & \left| \sum_{j=1}^p \left( \frac{\|\tilde{\epsilon}_j\|_n^2 - |\tilde{\omega}_j^0|^2}{|\tilde{\omega}_j^0|^2} \right) \left( \frac{|\tilde{\omega}_j^0|^2 - \hat{\omega}_j^2}{\hat{\omega}_j^2} \right) \right| \\ & \leq \left( \sum_{j=1}^p \left( \frac{\|\tilde{\epsilon}_j\|_n^2 - |\tilde{\omega}_j^0|^2}{|\tilde{\omega}_j^0|^2} \right)^2 \right)^{1/2} \left( \sum_{j=1}^p \left( \frac{|\tilde{\omega}_j^0|^2 - \hat{\omega}_j^2}{\hat{\omega}_j^2} \right)^2 \right)^{1/2} \\ & \leq 2\sqrt{(p + \tilde{s})\lambda_2^2} \left( \sum_{j=1}^p \left( \frac{|\tilde{\omega}_j^0|^2 - \hat{\omega}_j^2}{\hat{\omega}_j^2} \right)^2 \right)^{1/2} \leq \frac{(p + \tilde{s})\lambda_2^2}{\delta_2} + \delta_2 \sum_{j=1}^p \left( \frac{|\tilde{\omega}_j^0|^2 - \hat{\omega}_j^2}{\hat{\omega}_j^2} \right)^2. \end{aligned}$$

Invoking  $\text{trace}(\Theta_0 \Sigma_n) = \text{trace}(\tilde{\Theta}_0 \Sigma_n)$ , that is  $\sum_{j=1}^p \|\epsilon_j\|_n^2 / |\omega_j^0|^2 = \sum_{j=1}^p \|\tilde{\epsilon}_j\|_n^2 / |\tilde{\omega}_j^0|^2$ , and using that we are on  $\mathcal{T}_1$ , we see that

$$\begin{aligned} & \left( \frac{1}{K_0^2} - \delta_1 \right) \sum_{j=1}^p \|X(\hat{\beta}_j - \tilde{\beta}_j^0)\|_n^2 + \left( \frac{1}{2K_0^4 \sigma_0^4} - \delta_2 \right) \sum_{j=1}^p \left( \frac{\hat{\omega}_j^2 - |\tilde{\omega}_j^0|^2}{|\hat{\omega}_j|^2} \right)^2 \\ & \quad + \left( \lambda^2 - \frac{\lambda_1^2}{\delta_1} \right) \hat{s} \leq \lambda^2 s_0 + \frac{\lambda_2^2 (p + \tilde{s})}{\delta_2} + \frac{\lambda_1^2 \tilde{s}}{\delta_1}. \end{aligned}$$

□

### 7.3. Exploiting the beta-min condition.

**Lemma 7.2** *Let  $\tilde{s} = s_{\tilde{B}}$  be the number of edges of  $\tilde{B}_0$  and  $\hat{s} = s_{\hat{B}_0}$  be the number of edges of  $\hat{B}$ . Suppose that for some  $\tilde{\lambda}$ ,*

$$\|\hat{B} - \tilde{B}_0\|_F \leq \tilde{\lambda} \sqrt{\tilde{s}},$$

and that for some constant  $0 \leq \eta_1 < 1$  and  $0 < \eta_2^2 < 1 - \eta_1$

$$\#\{|\tilde{\beta}_{j,k}^0| \geq \tilde{\lambda}/\eta_2\} \geq (1 - \eta_1)\tilde{s}.$$

Then  $\hat{s} \geq (1 - \eta_1 - \eta_2^2)\tilde{s}$ .

**Proof.** Let

$$\mathcal{N} := \{(k, j) : |\tilde{\beta}_{k,j}^0| \geq \tilde{\lambda}/\eta_2\}, \quad \mathcal{M} := \{(k, j) : |\hat{\beta}_{k,j} - \tilde{\beta}_{k,j}^0| \geq \tilde{\lambda}/\eta_2\}.$$

Then for  $(k, j) \in \mathcal{N} \cap \mathcal{M}^c$  it holds that

$$|\hat{\beta}_{k,j}| \geq |\tilde{\beta}_{k,j}^0| - |\hat{\beta}_{k,j} - \tilde{\beta}_{k,j}^0| > 0,$$

so that  $\hat{s} \geq |\mathcal{N} \cap \mathcal{M}^c|$ . Since  $\|\hat{B} - \tilde{B}_0\|_F \leq \tilde{\lambda}\sqrt{\tilde{s}}$ , we must have

$$\sum_{(k,j) \in \mathcal{N} \cap \mathcal{M}} |\hat{\beta}_{k,j} - \tilde{\beta}_{k,j}^0|^2 \leq \sum_{(k,j)} |\hat{\beta}_{k,j} - \tilde{\beta}_{k,j}^0|^2 = \|\hat{B} - \tilde{B}_0\|_F^2 \leq \tilde{\lambda}^2 \tilde{s},$$

whereas

$$\sum_{(k,j) \in \mathcal{N} \cap \mathcal{M}} |\hat{\beta}_{k,j} - \tilde{\beta}_{k,j}^0|^2 \geq |\mathcal{N} \cap \mathcal{M}| \tilde{\lambda}^2 / \eta_2^2.$$

Hence  $|\mathcal{N} \cap \mathcal{M}| \leq \eta_2^2 \tilde{s}$ . This gives

$$|\mathcal{N} \cap \mathcal{M}^c| = |\mathcal{N}| - |\mathcal{N} \cap \mathcal{M}| \geq (1 - \eta_1) \tilde{s} - \eta_2^2 \tilde{s} = (1 - \eta_1 - \eta_2^2) \tilde{s}.$$

□

**Lemma 7.3** *Suppose that for some  $\delta_B > 0$ ,  $\delta_s > 0$ ,  $\lambda_0 > 0$  and  $\lambda$  one has*

$$\delta_B \|\hat{B} - \tilde{B}_0\|_F^2 + \lambda^2 \delta_s \hat{s} \leq \lambda^2 s_0 + \lambda_0^2 \tilde{s},$$

where  $\tilde{s} \geq s_0$ . Let  $\tilde{\lambda}^2 \delta_B \geq \lambda^2 + \lambda_0^2$  and assume that

$$\#\{|\tilde{\beta}_{j,k}^0| \geq \tilde{\lambda} / \eta_2\} \geq (1 - \eta_1) \tilde{s}.$$

Then

$$\delta_B \|\hat{B} - \tilde{B}_0\|_F^2 + \left( \lambda^2 \delta_s - \frac{\lambda_0^2}{1 - \eta_1^2 - \eta_2^2} \right) \hat{s} \leq \lambda^2 s_0,$$

and  $\hat{s} \geq (1 - \eta_1 - \eta_2^2) s_0$ .

**Proof.** Since  $\tilde{s} \geq s_0$ , we find that

$$\delta_B \|\hat{B} - \tilde{B}_0\|_F^2 \leq (\lambda^2 + \lambda_0^2) \tilde{s} \leq \delta_B \tilde{\lambda}^2 \tilde{s}.$$

This gives by Lemma 7.2 that  $\hat{s} \geq (1 - \eta_1 - \eta_2^2) \tilde{s}$ . But then

$$\delta_B \|\hat{B} - \tilde{B}_0\|_F^2 + \left( \lambda^2 \delta_s - \frac{\lambda_0^2}{1 - \eta_1 - \eta_2^2} \right) \hat{s} \leq \lambda^2 s_0.$$

□

7.4. *The sets  $\mathcal{T}_k$ ,  $k = 0, 1, 2, 3$ .*

7.4.1. *The set  $\mathcal{T}_1$ .*

**Lemma 7.4** *Let  $Z$  be a fixed  $n \times m$  matrix and  $\varepsilon_1, \dots, \varepsilon_n$  be independent  $\mathcal{N}(0, \sigma_0^2)$ -distributed random variables. Then for all  $t > 0$*

$$\mathbb{P}\left(\sup_{\|Z\beta\|_n \leq 1} |\varepsilon^T Z\beta|/n \geq \sigma_0(\sqrt{2m/n} + \sqrt{2t/n})\right) \leq \exp[-t].$$

**Proof.** Assume without loss of generality that  $Z^T Z/n = I$  and define  $V_k := \varepsilon^T Z_k / (\sigma_0 \sqrt{n})$ . Then  $V_1, \dots, V_p$  are independent and  $\mathcal{N}(0, 1)$ -distributed. It follows that for all  $N \in \{2, 3, \dots\}$ , that  $\mathbb{E}|V_k^2|^N = (2N)!/(2^N N!) \leq N!$ . But then by Bernstein's inequality (see Bennett (1962)), for all  $t > 0$ ,

$$(9) \quad \mathbb{P}\left(\sum_{k=1}^m (V_k^2 - \mathbb{E}V_k^2) \geq 2\sqrt{tm} + 2t\right) \leq \exp[-t].$$

Now use that  $\sum_{j=1}^m \mathbb{E}V_k^2 = m$ . We get

$$\mathbb{P}\left(\sum_{k=1}^m V_k^2 \geq m + 2\sqrt{tm} + 2t\right) \leq \exp[-t].$$

But  $m + 2\sqrt{tm} + 2t \leq (\sqrt{2m} + \sqrt{2t})^2$ . Furthermore,

$$\sup_{\|Z\beta\|_n \leq 1} |\varepsilon^T Z\beta|/n = \frac{\sigma_0}{\sqrt{n}} \sqrt{\sum_{k=1}^m V_k^2/n}.$$

□

We are dealing now with the problem of uniformly controlling over all permutations  $\pi$ . We consider the local structure at each node of a DAG  $(\tilde{B}_0(\pi), \tilde{\Omega}_0(\pi))$  with  $\tilde{B}_0(\pi) =: (\tilde{\beta}_{k,j}^0(\pi))$ . Let  $\tilde{S}_j(\pi)$  be the set of incoming edges at node  $j$ . Given  $\tilde{S}_j(\pi)$ , the vector  $X\tilde{\beta}_j^0(\pi)$  is the projection in  $L_2(P)$  of  $X_j$  on the linear space spanned by  $\{X_k\}_{k \in \tilde{S}_j(\pi)}$ . Moreover,  $\tilde{\varepsilon}_j(\pi)$  is the anti-projection  $\tilde{\varepsilon}_j(\pi) = X_j - X\tilde{\beta}_j^0(\pi)$ . In other words (for  $j$  fixed) if the parents  $\tilde{S}_j(\pi)$  at node  $j$  are given, then the coefficients  $\tilde{\beta}_{k,j}^0(\pi)$  and noise term  $\tilde{\varepsilon}_j(\pi)$  are given as well. Also, the set of variables  $X_k$  that are independent of  $\tilde{\varepsilon}_j$  is then given. Recall that  $\tilde{B}_j(\pi) := \{\beta_j : X_k \perp \tilde{\varepsilon}_j(\pi), \forall \beta_{k,j} \neq 0\}$ . Thus, for each fixed  $j$ , if  $\tilde{S}_j(\pi)$  is given then the local situation  $(\tilde{\varepsilon}_j(\pi), \tilde{\beta}_j^0(\pi), \tilde{B}_j(\pi))$  at node  $j$  is given.

Let  $\Pi_j(m)$  be the collection of all permutations giving DAGs  $(\tilde{B}_0(\pi), \tilde{\Omega}_0(\pi))$  with edges  $(\tilde{S}_1(\pi), \dots, \tilde{S}_p(\pi))$  with  $|\tilde{S}_j(\pi)| = m$ . If for some  $m \in \{0, 1, \dots, p\}$ , we know that  $\pi \in \Pi_j(m)$ , that is, we know that node  $j$  has  $m$  parents, then there are at most  $\binom{p}{m}$  possibilities for the local situation at node  $j$ .

**Theorem 7.1** *Assume Condition 3.1. Then for all  $t > 0$ ,*

$$\begin{aligned} \mathbb{P} \left( \max_{\pi} \sup_{B \in \tilde{\mathcal{B}}(\pi)} 2 \sum_{j=1}^p |\tilde{\epsilon}_j^T(\pi)(X(\beta_j - \tilde{\beta}_j^0(\pi))|/n - \delta_1 \sum_{j=1}^p \|X(\beta_j - \tilde{\beta}_j^0(\pi))\|_n^2 \right. \\ \left. \geq \frac{4\sigma_0^2(s_B + \tilde{s}(\pi))}{n\delta_1} + \frac{\sigma_0^2(t + 2 \log p)(s_B + \tilde{s}(\pi))}{n\delta_1} \right) \\ \leq \exp[-t]. \end{aligned}$$

**Proof.** Let  $A_j(\pi)$  be the event

$$\begin{aligned} A_j(\pi) &:= \left\{ \exists \beta_j \in \tilde{\mathcal{B}}_j(\pi) : \sup_{\|X(\beta_j - \tilde{\beta}_j^0(\pi))\|_n \leq 1} |\tilde{\epsilon}_j^T(\pi)(X(\beta_j - \tilde{\beta}_j^0(\pi))|/n \right. \\ &\left. \geq \sigma_0 \left( \sqrt{\frac{2(s_{\beta_j} + \tilde{s}_j(\pi))}{n}} + \sqrt{\frac{2(t + \tilde{s}_j(\pi) \log p + 2 \log p)}{n}} \right) \right\}. \end{aligned}$$

Then by Lemma 7.4, for all  $t > 0$ ,  $\pi$  and  $j$

$$\mathbb{P}(A_j(\pi)) \leq \exp[-(t + \tilde{s}_j(\pi) \log p + 2 \log p)].$$

We now let  $\pi$  vary over all permutations such that  $|\tilde{S}_j(\pi)| = m$ . We then get

$$\mathbb{P} \left( \bigcup_{\pi \in \Pi_j(m)} A_j(\pi) \right) \leq \binom{p}{m} \exp[-(t + m \log p + 2 \log p)] \leq \exp[-(t + 2 \log p)].$$

Next, we let  $\pi$  vary over all permutations. We get

$$\begin{aligned} \mathbb{P} \left( \bigcup_{\pi} A_j(\pi) \right) &\leq \sum_{m=1}^p \max_{1 \leq m \leq p} \mathbb{P} \left( \bigcup_{\pi \in \Pi_j(m)} A_j(\pi) \right) \\ &\leq p \exp[-(t + 2 \log p)] \leq \exp[-(t + \log p)]. \end{aligned}$$

Finally

$$\mathbb{P}(\bigcup_{j=1}^p \bigcup_{\pi} A_j(\pi)) \leq p \max_j \mathbb{P}(\bigcup_{\pi} A_j(\pi)) \leq p \exp[-(t + \log p)] \leq \exp[-t].$$

Now, we use that for all  $\delta_1 > 0$ ,

$$\begin{aligned} & 2\sigma_0 \sum_{j=1}^p \left( \sqrt{\frac{2(s_j + \tilde{s}_j)}{n}} + \sqrt{\frac{2(t + \tilde{s}_j + 2 \log p)}{n}} \right) \|X(\beta_j - \tilde{\beta}_j^0)\|_n \\ & \leq \delta_1 \sum_{j=1}^p \|X(\beta_j - \tilde{\beta}_j^0)\|_n^2 + \frac{4\sigma_0^2(s + \tilde{s})}{n\delta_1} + \frac{4\sigma_0^2(t + 2 \log p)(s + \tilde{s})}{n\delta_1}, \end{aligned}$$

where  $s = \sum_{j=1}^p s_j$ ,  $\tilde{s} = \sum_{j=1}^p \tilde{s}_j$ .

□

#### 7.4.2. The set $\mathcal{I}_2$ .

**Theorem 7.2** *Assume Condition 3.4. Then for all  $t > 0$ ,*

$$\begin{aligned} & \mathbb{P} \left( \exists \pi : \sum_{j=1}^p \left( \frac{\|\tilde{\epsilon}_j(\pi)\|_n^2 - |\tilde{\omega}_j^0(\pi)|^2}{|\tilde{\omega}_j^0(\pi)|^2} \right)^2 \right. \\ & \geq 8 \left( \frac{pt + (1 + 8\tilde{\alpha})\tilde{s}(\pi) \log p + 2p \log p}{n} \right) + 8 \left( \frac{4p(t^2 + \log^2 p)}{n^2} \right) \Bigg) \\ & \leq 2 \exp[-t]. \end{aligned}$$

**Proof.** Define

$$Z_j(\pi) := \frac{\|\tilde{\epsilon}_j(\pi)\|_n^2 - |\tilde{\omega}_j^0(\pi)|^2}{|\tilde{\omega}_j^0(\pi)|^2}.$$

Using the same argument as in (9), we see that for each  $\pi$ , and for all  $t > 0$ .

$$\mathbb{P} \left( |Z_j(\pi)| \geq 2 \left( \sqrt{\frac{t}{n}} + \frac{t}{n} \right) \right) \leq 2 \exp[-t].$$

Define  $\mathbf{Z}_j(\pi) := |Z_j(\pi)|/2a_j(\pi)$ , where

$$\begin{aligned} a_j(\pi) &= \left( \sqrt{\frac{t + \tilde{s}_j(\pi) \log p + \log(1 + p) + \log p}{n}} \right. \\ & \quad \left. + \frac{t + \tilde{s}_j(\pi) \log p + \log(1 + p) + \log p}{n} \right). \end{aligned}$$

It follows that

$$\mathbb{P} \left( \max_{1 \leq j \leq p} \max_{0 \leq m \leq p} \max_{\pi \in \Pi_j(m)} \mathbf{Z}_j(\pi) \geq 1 \right)$$

$$\leq 2p(p+1) \binom{p}{m} \exp[-(t + m \log p + \log(1+p) + \log p)] \leq 2 \exp[-t].$$

Invoking  $\log(1+p) \leq 2 \log p$ , we see that with probability at least  $1 - 2 \exp[-t]$ , it holds that for all permutations  $\pi$  and all  $j$ ,

$$|Z_j(\pi)| \leq 2 \sqrt{\frac{t + \tilde{s}_j(\pi) \log p + 2 \log p}{n}} + \frac{t + \tilde{s}_j(\pi) \log p + 2 \log p}{n},$$

which implies

$$\begin{aligned} \sum_{j=1}^p |Z_j(\pi)|^2 &\leq 4 \sum_{j=1}^p \left( \sqrt{\frac{t + \tilde{s}_j(\pi) \log p + 2 \log p}{n}} + \frac{t + \tilde{s}_j(\pi) \log p + 2 \log p}{n} \right)^2 \\ &\leq 8 \left( \frac{pt + \tilde{s}(\pi) \log p + 2p \log p}{n} \right) + 8 \left( \frac{4pt^2 + 8 \sum_{j=1}^p \tilde{s}_j^2(\pi) \log^2 p + 4p \log^2 p}{n^2} \right). \end{aligned}$$

Next, we insert that for all  $j$ ,  $\tilde{s}_j(\pi) \leq \tilde{\alpha}n/(\log p)$ , to find

$$\sum_{j=1}^p \tilde{s}_j^2(\pi) \log^2 p \leq \sum_{j=1}^n (\tilde{\alpha}n/(\log p)) \tilde{s}_j(\pi) \log^2 p = \tilde{\alpha} \tilde{s}(\pi) n \log p.$$

We then arrive at

$$\sum_{j=1}^p |Z_j(\pi)|^2 \leq 8 \left( \frac{pt + (1 + 8\tilde{\alpha}) \tilde{s}(\pi) \log p + 2p \log p}{n} \right) + 8 \left( \frac{4p(t^2 + \log^2 p)}{n^2} \right).$$

□

### 7.4.3. The sets $\mathcal{T}_3$ and $\mathcal{T}_0$ .

**Theorem 7.3** *Assume Condition 3.1 and Condition 3.2. For all  $t > 0$ , with probability at least  $1 - 2 \exp[-t]$ ,*

$$\|X\beta\|_n \geq \left[ 3\Lambda_{\min}/4 - \sqrt{\frac{2(t + \log p)}{n}} - 3\sigma_0 \sqrt{\frac{s_\beta \log p}{n}} \right] \|\beta\|_2,$$

uniformly in all  $\beta \in \mathbb{R}^p$ .

**Proof.**

We follow here the arguments used in Raskutti, Wainwright and Yu (2010), which we slightly adjust to the style of the present paper. They show that for  $\delta'_3 = 1/4$  (in fact for  $\delta'_3 = o(1)$  as  $n \rightarrow \infty$ ), and for all  $r > 0$

$$\mathbb{E} \inf_{\|\beta\|_1 \leq r, \|X\beta\|=1} \|X\beta\|_n \geq 1 - \delta'_3 - 3\sigma_0 \sqrt{\frac{\log p}{n}} r.$$

Hence, for all  $1 \leq m \leq p$ ,

$$\mathbf{E} \inf_{s_\beta \leq m, \|\beta\|_2=1} \|X\beta\|_n \geq (1 - \delta'_3)\Lambda_{\min} - 3\sigma_0 \sqrt{\frac{m \log p}{n}}.$$

Apply the concentration inequality given in Massart (2003) to find that for all  $t > 0$ ,

$$\mathbf{P} \left( \left[ \mathbf{E} \inf_{s_\beta \leq m, \|\beta\|_2=1} \|X\beta\|_n \right] - \left[ \inf_{s_\beta \leq m, \|\beta\|_2=1} \|X\beta\|_n \right] \geq \sqrt{\frac{2t}{n}} \right) \leq 2 \exp[-t].$$

Thus

$$\mathbf{P} \left( \left[ (1 - \delta'_3)\Lambda_{\min} - 3\sigma_0 \sqrt{\frac{m \log p}{n}} \right] - \left[ \inf_{s_\beta \leq m, \|\beta\|_2=1} \|X\beta\|_n \right] \geq \sqrt{\frac{2t}{n}} \right) \leq 2 \exp[-t],$$

and hence

$$\begin{aligned} \mathbf{P} \left( \exists \beta : \left[ (1 - \delta'_3)\Lambda_{\min} - 3\sigma_0 \sqrt{\frac{s_\beta \log p}{n}} \right] \|\beta\|_2 - \|X\beta\|_n \right. \\ \left. \geq \sqrt{\frac{2(t + \log p)}{n}} \|\beta\|_2 \right) \leq 2 \exp[-t]. \end{aligned}$$

□

**Lemma 7.5** *Assume Conditions 3.1, 3.2, 3.3, and 3.4 and that*

$$1/K_0 := 3\Lambda_{\min}/4 - \sqrt{\frac{2(t + \log p)}{n}} - 3\sigma_0 \sqrt{\alpha + \tilde{\alpha}} > 0.$$

Let for some  $t > 0$ .

$$\tilde{\mathcal{T}}_3 := \left\{ \|X\beta\|_n \leq \left[ 3\Lambda_{\min}/4 - \sqrt{\frac{2(t + \log p)}{n}} - 3\sigma_0 \sqrt{\frac{s_\beta \log p}{n}} \right] \|\beta\|_2, \forall \beta \right\}.$$

Then  $\mathbf{P}(\tilde{\mathcal{T}}_3) \geq 1 - 2 \exp[-t]$  and one has on  $\tilde{\mathcal{T}}_3$ , for all  $B = (\beta_1, \dots, \beta_p) \in \mathcal{B}$  and all  $\pi$  and all  $j$ ,

$$(10) \quad \|X(\beta_j - \tilde{\beta}_j^0(\pi))\|_n \geq \|\beta_j - \tilde{\beta}_j^0(\pi)\|_2 / K_0^2.$$

Moreover, on  $\tilde{\mathcal{T}}_3$ , also  $\hat{\omega}_j^2 \geq 1/K_0^2$  for all  $j$ .

**Proof.** Theorem 7.3 states that  $\mathbf{P}(\tilde{\mathcal{T}}_3) \geq 1 - 2 \exp[-t]$ . Result (10) follows immediately, since  $s_{\beta_j} + s_{\tilde{\beta}_j^0} \leq (\alpha + \tilde{\alpha})n / \log p$ . For the last result, we define  $\hat{\beta}_{k,j}^- := -\hat{\beta}_{k,j}$  for  $k \neq j$  and  $\hat{\beta}_{j,j}^- = 1$ . Then on  $\tilde{\mathcal{T}}_3$ ,

$$\hat{\omega}_j^2 = \|X\hat{\beta}_j^-\|_n^2 \geq \|\hat{\beta}_j^-\|_n / K_0^2 \geq 1/K_0^2.$$

□

7.5. *Collecting the results.*

**Lemma 7.6** *Define  $(\tilde{B}_0, \tilde{\Omega}_0) := (B_0(\hat{\pi}), \Omega_0(\hat{\pi}))$ . Assume Conditions 3.1, 3.2, 3.3, 3.4, and 3.5. Suppose we are on  $\cap_{k=0}^3 \mathcal{T}_k$  with  $0 < \delta_1 < 1/K_0^2$  and  $0 < \delta_2 < 1/(2K_0^4\sigma_0^4)$  and  $\delta_3 - \lambda_3\sqrt{\alpha} + \tilde{\alpha}\sqrt{n/\log p} \geq 1/K_0 > 0$ . Take the tuning parameter  $\lambda^2 > \lambda_1^2/\delta_1 + \lambda_2^2/\delta_2$ . Let*

$$\begin{aligned} \delta_B &\leq \frac{1}{K_0^2} \left( \frac{1}{K_0^2} - \delta_1 \right), \quad \delta_W \leq \frac{1}{K_0^2} \left( \frac{1}{2K_0^4\sigma_0^4} - \delta_2 \right), \\ \delta_s &\leq \left( 1 - \frac{\lambda_1^2}{\lambda^2\delta_1} - \frac{\lambda_2^2}{\lambda^2\delta_2} \right), \quad \lambda_0^2 := \frac{(p/s_0 + 1)\lambda_2^2}{\delta_2} + \frac{\lambda_1^2}{\delta_1}. \end{aligned}$$

Let  $\tilde{\lambda}^2\delta_B := \lambda^2 + \lambda_0^2$ , and  $\eta_2^2 := \eta_0^2\tilde{\lambda}^2n/\log p = \eta_0^2(\lambda^2 + \lambda_0^2)(n/\log p)/\delta_B^2$ . Assume

$$\left( \lambda^2\delta_s - \frac{\lambda_0^2}{1 - \eta_1 - \eta_2^2} \right) := \lambda^2\delta_\eta > 0.$$

Then

$$\delta_B \|\hat{B} - \tilde{B}_0\|_F^2 + \delta_W \|\hat{\Omega} - \tilde{\Omega}_0\|_F^2 + \lambda^2\delta_\eta \hat{s} \leq \lambda^2 s_0,$$

and  $\hat{s} \geq (1 - \eta_1 - \eta_2^2)\tilde{s} \geq (1 - \eta_1 - \eta_2^2)s_0$ .

**Proof.** This follows from Lemma 7.1 and Lemma 7.3.  $\square$

**Lemma 7.7** *Assume Conditions 3.1, 3.2, 3.3 and 3.4, with*

$$3\Lambda_{\min}/4 - \sqrt{\frac{2(t + \log p)}{n}} - 3\sigma_0\sqrt{\alpha} + \tilde{\alpha} \geq 1/K_0 > 0.$$

Take

$$\begin{aligned} \lambda_1^2 &:= \frac{4\sigma_0^2(1 + t + 2\log p) \log p}{\log p \cdot n}, \\ \lambda_2^2 &:= 4 \left( 3 + 8\tilde{\alpha} + \frac{(t + 2\log p)}{\log p} + \frac{4(t^2 + \log^2 p)}{n \log p} \right) \frac{\log p}{n}, \\ \lambda_3^2 &:= 9\sigma_0^2 \frac{\log p}{n}, \quad \delta_3 := \frac{3}{4}\Lambda_{\min} - \sqrt{\frac{2(t + \log p)}{n}}. \end{aligned}$$

Then

$$\mathbb{P}(\cap_{k=0}^3 \mathcal{T}_k) \geq 1 - 4 \exp[-t].$$

**Proof.** This follows from combining Theorem 7.1, Theorem 7.2 and Lemma 7.5.  $\square$

**Theorem 7.4** *Assume Conditions 3.1, 3.2, 3.3, 3.4 and 3.5. Let us take  $t = \log p$ , giving  $\alpha_0 = 4/p$  (suppose  $p$  is large). Take  $n$  sufficiently large, and  $\log p/n$  bounded. Let*

$$c_1 := 96, \quad c_2 := 3840, \quad c = 4 \left( \frac{(p/s_0 + 1)c_2\sigma_0^2}{\Lambda_{\min}^4} + \frac{c_1\sigma_0^2}{\Lambda_{\min}^2} \right) + 2 \left( \frac{c_1\sigma_0^2}{\Lambda_{\min}^2} + \frac{c_2\sigma_0^2}{\Lambda_{\min}^4} \right).$$

*Some possible choices for the constants are*

$$\alpha = \tilde{\alpha} = \frac{\Lambda_{\min}^2}{288\sigma_0^2}, \quad K_0 := \frac{2}{\Lambda_{\min}},$$

$$\delta_1 := \frac{\Lambda_{\min}^2}{8}, \quad \delta_2 := \frac{\Lambda_{\min}^4}{64\sigma_0^4}, \quad \delta_3 := \frac{\Lambda_{\min}}{2}.$$

$$\lambda^2 := c \frac{\log p}{n}, \quad \lambda_1^2 = 12\sigma_0^2 \frac{\log p}{n}, \quad \lambda_2^2 = 60 \frac{\log p}{n}, \quad \lambda_3^2 = 9\sigma_0^2 \frac{\log p}{n}.$$

*Then*

$$\delta_B = \frac{\Lambda_{\min}^4}{32}, \quad \delta_W = \frac{\Lambda_{\min}^6}{256\sigma_0^4}, \quad \delta_s = \left( 1 - \frac{c_1\sigma_0^2}{c\Lambda_{\min}^2} - \frac{c_2\sigma_0^4}{c\Lambda_{\min}^4} \right).$$

*We let*

$$\lambda_0^2 := \left( \frac{(p/s_0 + 1)c_2\sigma_0^4}{\Lambda_{\min}^4} + \frac{c_1\sigma_0^2}{\Lambda_{\min}^2} \right) \frac{\log p}{n}.$$

$$\tilde{\lambda}^2 = \left( c + \frac{(p/s_0 + 1)c_2\sigma_0^4}{\Lambda_{\min}^4} + \frac{c_1\sigma_0^2}{\Lambda_{\min}^2} \right) \frac{32}{\Lambda_{\min}^4} \frac{\log p}{n},$$

*and*

$$\delta_\eta = \frac{1}{2}, \quad \eta_1 = 0, \quad 2\eta_0^2 = \left( c + \frac{(p/s_0 + 1)c_2\sigma_0^4}{\Lambda_{\min}^4} + \frac{c_1\sigma_0^2}{\Lambda_{\min}^2} \right)^{-1},$$

$$\eta_2^2 = \eta_0^2 \left( c + \frac{(p/s_0 + 1)c_2\sigma_0^4}{\Lambda_{\min}^4} + \frac{c_1\sigma_0^2}{\Lambda_{\min}^2} \right) \frac{32}{\Lambda_{\min}^4}.$$

**Proof.** This follows from using some bounds and exact choices in Lemma 7.6 and Lemma 7.7. In particular, with  $\lambda^2 = c \log p/n$ , we take  $\tilde{\lambda}^2 = (\lambda^2 + \lambda_0^2)/\delta_B$ . With  $\eta_1 = 0$  and  $\delta_\eta = 1/2$ , the equation

$$\left( \lambda^2 \delta_s - \frac{\lambda_0^2}{1 - \eta_2^2} \right) = \frac{\lambda^2}{2}$$

gives

$$c \left( 1 - \frac{c_1\sigma_0^2}{c\Lambda_{\min}^2} - \frac{c_2\sigma_0^4}{c\Lambda_{\min}^4} \right)$$

$$-\left(\frac{(p/s_0 + 1)c_2\sigma_0^4}{\Lambda_{\min}^4} + \frac{c_1\sigma_0^2}{\Lambda_{\min}^2}\right) \left(1 - \eta_1 - \eta_0^2 \left(c + \frac{(p/s_0 + 1)c_2\sigma_0^4}{\Lambda_{\min}^4} + \frac{c_1\sigma_0^2}{\Lambda_{\min}^2}\right)\right)^{-1} = \frac{c}{2}.$$

With

$$2\eta_0^2 = \left(c + \frac{c_2\sigma_0^4}{\Lambda_{\min}^4} + \frac{c_1\sigma_0^2}{\Lambda_{\min}^2}\right)^{-1}$$

we have to solve for  $c$

$$c \left(1 - \frac{c_1\sigma_0^2}{c\Lambda_{\min}^2} - \frac{c_2\sigma_0^4}{c\Lambda_{\min}^4}\right) - 2 \left(\frac{(p/s_0 + 1)c_2\sigma_0^4}{\Lambda_{\min}^4} + \frac{c_1\sigma_0^2}{\Lambda_{\min}^2}\right) = \frac{c}{2}.$$

This yields

$$c = 4 \left(\frac{(p/s_0 + 1)c_2\sigma_0^4}{\Lambda_{\min}^4} + \frac{c_1\sigma_0^2}{\Lambda_{\min}^2}\right) + 2 \left(\frac{c_1\sigma_0^2}{\Lambda_{\min}^2} + \frac{c_2\sigma_0^4}{\Lambda_{\min}^4}\right).$$

□

**7.6. Proof of Theorem 5.1.** We investigate what happens on the set  $\cap_{k=1}^3 \mathcal{I}_3$  defined in Subsection 7.2. The results in Subsection 7.4 say that  $\cap_{k=1}^3 \mathcal{I}_k$  has probability at least  $4\exp[-t]$  for a proper choice of the constants and parameters involved. Theorem 5.1 then follows directly.

**Lemma 7.8** *Assume Condition 3.1, Condition 5.1 and Condition 5.2. Suppose we are on  $\cap_{k=1}^3 \mathcal{I}_k$ , with*

$$\lambda^2 > \lambda_1^2/\delta_1,$$

$$\delta_3 - \lambda_3\sqrt{2\alpha_*} \geq \frac{1}{K_0} > 0,$$

and

$$(11) \quad \left(\frac{1}{2\sigma_0^4} - \delta_2\right) \geq \eta_\omega \left(\frac{2\lambda_1^2}{\delta_2} + \frac{\lambda_1^2}{\delta_1} + \lambda^2\right) \frac{\alpha_* n}{\log p}.$$

Then  $\hat{\pi} = \pi_0$  and

$$\left(\frac{1 - \delta_1}{K_0^2}\right) \|\hat{B} - B_0\|_F^2 + \left(\lambda^2 - \frac{\lambda_1^2}{\delta_1}\right) \hat{s} \leq \left(\lambda^2 + \frac{\lambda_1^2}{\delta_1}\right) s_0.$$

**Proof.** We have

$$\sum_{j=1}^p \|X_j - X\hat{\beta}_j\|_n^2 + \lambda^2 \hat{s} \leq \sum_{j=1}^p \|\epsilon_j\|_n^2 + \lambda^2 s_0 = \sum_{j=1}^p \frac{\|\tilde{\epsilon}_j\|_n^2}{|\tilde{\omega}_j^0|^2}.$$

So we find

$$\sum_{j=1}^p \|X(\hat{\beta}_j - \tilde{\beta}_j^0)\|_n^2 + \lambda^2 \hat{s} \leq 2 \sum_{j=1}^p \tilde{\epsilon}_j^T X(\hat{\beta}_j - \tilde{\beta}_j^0)/n + \sum_{j=1}^p \|\tilde{\epsilon}_j\|_n^2 \left( \frac{1}{|\tilde{\omega}_j^0|^2} - 1 \right) + \lambda^2 s_0.$$

We have

$$\sum_{j=1}^p \|\tilde{\epsilon}_j\|_n^2 \left( \frac{1}{|\tilde{\omega}_j^0|^2} - 1 \right) = \sum_{j=1}^p \frac{\|\tilde{\epsilon}_j\|_n^2 - |\tilde{\omega}_j^0|^2}{|\tilde{\omega}_j^0|^2} (|\tilde{\omega}_j^0|^2 - 1) + \sum_{j=1}^p (1 - |\tilde{\omega}_j^0|^2).$$

We know that

$$\log(\det(\Sigma_0)) = \sum_{j=1}^p \log |\tilde{\omega}_j^0|^2 = \sum_{j=1}^p \log |\omega_j^0|^2 = 0,$$

since  $|\omega_j^0|^2 = 1$  for all  $j$ . Moreover

$$\log(1+x) \leq x - \frac{1}{2(1+c)^2} x^2, \quad -1 < x \leq c.$$

So, since  $|\tilde{\omega}_j^0|^2 \leq \sigma_0^2$ ,

$$\log |\tilde{\omega}_j^0|^2 \leq (|\tilde{\omega}_j^0|^2 - 1) - \frac{1}{2\sigma_0^4} (|\tilde{\omega}_j^0|^2 - 1)^2.$$

Hence

$$0 \leq \sum_{j=1}^p (|\tilde{\omega}_j^0|^2 - 1) - \frac{1}{2\sigma_0^4} \sum_{j=1}^p (|\tilde{\omega}_j^0|^2 - 1)^2.$$

This gives

$$\sum_{j=1}^p (1 - |\tilde{\omega}_j^0|^2) \leq -\frac{1}{2\sigma_0^4} \sum_{j=1}^p (|\tilde{\omega}_j^0|^2 - 1)^2.$$

Therefore

$$\begin{aligned} & \sum_{j=1}^p \|X(\hat{\beta}_j - \tilde{\beta}_j^0)\|_n^2 + \frac{1}{2\sigma_0^4} \sum_{j=1}^p (|\tilde{\omega}_j^0|^2 - 1)^2 + \lambda^2 \hat{s} \\ & \leq 2 \sum_{j=1}^p \tilde{\epsilon}_j^T X(\hat{\beta}_j - \tilde{\beta}_j^0)/n + \lambda^2 s_0 + \sum_{j=1}^p \frac{\|\tilde{\epsilon}_j\|_n^2 - |\tilde{\omega}_j^0|^2}{|\tilde{\omega}_j^0|^2} (|\tilde{\omega}_j^0|^2 - 1) \\ & \leq 2 \sum_{j=1}^p \tilde{\epsilon}_j^T X(\hat{\beta}_j - \tilde{\beta}_j^0)/n + \lambda^2 s_0 + \frac{\lambda_2^2 p}{\delta_2} + \delta_2 \sum_{j=1}^p (|\tilde{\omega}_j^0|^2 - 1)^2, \end{aligned}$$

where we invoked that we are on the set  $\mathcal{T}_2$ . We find

$$\begin{aligned} & \sum_{j=1}^p \|X(\hat{\beta}_j - \tilde{\beta}_j^0)\|_n^2 + \left(\frac{1}{2\sigma_0^4} - \delta_2\right) \sum_{j=1}^p (|\tilde{\omega}_j^0|^2 - 1)^2 + \lambda^2 \hat{s} \\ & \leq 2 \sum_{j=1}^p \tilde{\epsilon}_j^T X(\hat{\beta}_j - \tilde{\beta}_j^0)/n + \lambda^2 s_0 + \frac{\lambda_2^2 p}{\delta_2}. \end{aligned}$$

This gives in a next step, using that we are on  $\mathcal{T}_1$ ,

$$\begin{aligned} & (1 - \delta_1) \sum_{j=1}^p \|X(\hat{\beta}_j - \tilde{\beta}_j^0)\|_n^2 + \left(\frac{1}{2\sigma_0^4} - \delta_2\right) \sum_{j=1}^p (|\tilde{\omega}_j^0|^2 - 1)^2 + \lambda^2 \hat{s} \\ & \leq \frac{\lambda_2^2(p + \tilde{s})}{\delta_2} + \frac{\lambda_1^2(\tilde{s} + \hat{s})}{\delta_1} + \lambda^2 s_0. \end{aligned}$$

Hence, using that we are on  $\mathcal{T}_3$  and invoking Condition 5.2,

$$\begin{aligned} & \frac{1}{K_0^2} (1 - \delta_1) \|\hat{B} - B_0\|_F^2 + \left(\frac{1}{2\sigma_0^4} - \delta_2\right) \sum_{j=1}^p (|\tilde{\omega}_j^0|^2 - 1)^2 + \left(\lambda^2 - \frac{\lambda_1^2}{\delta_1}\right) \hat{s} \\ & \leq \frac{\lambda_2^2(p + \tilde{s})}{\delta_2} + \frac{\lambda_1^2 \tilde{s}}{\delta_1} + \lambda^2 s_0 \\ & \leq \left(\frac{2\lambda_2^2}{\delta_2} + \frac{\lambda_1^2}{\delta_1} + \lambda^2\right) p^2, \end{aligned}$$

where we use that  $\tilde{s} \leq p^2$  and  $s_0 \leq p^2$  (and also  $p \leq p^2$ ). Since (using again Condition 5.2)  $p \log p/n \leq \alpha_*$ , and

$$\sum_{j=1}^p \left(|\tilde{\omega}_j^0|^2 - 1\right)^2 > p/\eta_\omega \text{ if } \hat{\pi} \neq \pi_0,$$

find that if  $\hat{\pi} \neq \pi_0$ ,

$$\left(\frac{1}{2\sigma_0^4} - \delta_2\right) \frac{p}{\eta_\omega} < \left(\frac{2\lambda_2^2}{\delta_2} + \frac{\lambda_1^2}{\delta_1} + \lambda^2\right) \alpha_* \frac{n}{\log p},$$

which is in contradiction with Condition 5.1 and the further condition (11) imposed in this lemma. So we must have  $\hat{\pi} = \pi_0$ , and thus  $\tilde{\omega}_j^0 = 1$  for all  $j$ . The result now follows from restarting the proof with  $\tilde{\omega}_j^0 = 1$  for all  $j$  plugged in.

□

## References.

- ANDERSSON, S. A., MADIGAN, D. and PERLMAN, M. D. (1997). A characterization of Markov equivalence classes for acyclic digraphs. *The Annals of Statistics* **25** 505–541.
- BENNETT, G. (1962). Probability inequalities for sums of independent random variables. *Journal of the American Statistical Association* **57** 33–45.
- BÜHLMANN, P. and VAN DE GEER, S. (2011). *Statistics for High-Dimensional Data: Methods, Theory and Applications*. Springer.
- CHICKERING, D. M. (2002). Optimal structure identification with greedy search. *Journal of Machine Learning Research* **3** 507–554.
- HAUSER, A. and BÜHLMANN, P. (2012). Characterization and greedy learning of interventional Markov equivalence classes of directed acyclic graphs. *Journal of Machine Learning Research* **13** 2409–2464.
- KALISCH, M. and BÜHLMANN, P. (2007). Estimating high-dimensional directed acyclic graphs with the PC-algorithm. *Journal of Machine Learning Research* **8** 613–636.
- LAURITZEN, S. (1996). *Graphical Models*. Oxford University Press.
- LOH, P. L. and WAINWRIGHT, M. J. (2012). High-dimensional regression with noisy and missing data: Provable guarantees with non-convexity. *The Annals of Statistics* **40** 1637–1664.
- MAATHUIS, M. H., KALISCH, M. and BÜHLMANN, P. (2009). Estimating high-dimensional intervention effects from observational data. *The Annals of Statistics* **37** 3133–3164.
- MAATHUIS, M. H., COLOMBO, D., KALISCH, M. and BÜHLMANN, P. (2010). Predicting causal effects in large-scale systems from observational data. *Nature Methods* **7** 247–248.
- MASSART, P. (2003). Concentration inequalities and model selection. Lectures from the 33rd Summer School on Probability Theory held in Saint-Flour, July 6–23, 2003. *Lecture Notes in Mathematics*.
- PEARL, J. (2000). *Causality: Models, Reasoning and Inference*. Cambridge University Press.
- PETERS, J. and BÜHLMANN, P. (2012). Identifiability of Gaussian structural equation models with same error variances. arXiv:1205.2536.
- RASKUTTI, G., WAINWRIGHT, M. J. and YU, B. (2010). Restricted eigenvalue properties for correlated Gaussian designs. *Journal of Machine Learning Research* **11** 2241–2259.
- ROBINS, J. M., SCHEINES, R., SPRITES, P. and WASSERMAN, L. (2003). Uniform consistency in causal inference. *Biometrika* **90** 491–515.
- SHOJAIE, A. and MICHAELIDIS, G. (2010). Penalized likelihood methods for estimation of sparse high-dimensional directed acyclic graphs. *Biometrika* **97** 519–538.
- SILANDER, T. and MYLLYMÄKI, P. (2006). A simple approach for finding the globally optimal Bayesian network structure. In *UAI* 445–452.
- SPIRITES, P., GLYMOUR, C. and SCHEINES, R. (2000). *Causation, Prediction, and Search*, Second ed. MIT Press.
- UHLER, C., RASKUTTI, G., BÜHLMANN, P. and YU, B. (2012). Geometry of faithfulness assumption in causal inference. arXiv:1205.5473.
- VAN DE GEER, S., BÜHLMANN, P. and ZHOU, S. (2011). The adaptive and thresholded Lasso for potentially misspecified models (and a lower bound for the Lasso). *Electronic Journal of Statistics* **5** 688–749.
- ZHANG, J. and SPIRITES, P. (2003). Strong Faithfulness and Uniform Consistency in Causal Inference. In *UAI* 632–639.
- ZOU, H. (2006). The adaptive Lasso and its oracle properties. *Journal of the American Statistical Association* **101** 1418–1429.

SEMINAR FÜR STATISTIK  
ETH ZÜRICH  
E-MAIL: geer@stat.math.ethz.ch

SEMINAR FÜR STATISTIK  
ETH ZÜRICH  
E-MAIL: buhlmann@stat.math.ethz.ch