# Inference in high-dimensional graphical models

**Jana Janková**        **Sara van de Geer**

*Seminar for Statistics*
*ETH Zürich*

### Abstract

We provide a selected overview of methodology and theory for estimation and inference on the edge weights in high-dimensional *directed* and *undirected* Gaussian graphical models. For undirected graphical models, two main explicit constructions are provided: one based on a global method that maximizes the joint likelihood (the graphical Lasso) and one based on a local (nodewise) method that sequentially applies the Lasso to estimate the neighbourhood of each node. The estimators lead to confidence intervals for edge weights and recovery of the edge structure. We evaluate their empirical performance in an extensive simulation study. The theoretical guarantees for the methods are achieved under a sparsity condition relative to the sample size and regularity conditions. For directed acyclic graphs, we apply similar ideas to construct confidence intervals for edge weights, when the directed acyclic graph is identifiable.

## 1 Undirected graphical models

### 1.1 Introduction

Undirected graphical models, also known as Markov random fields, have become a popular tool for representing network structure of high-dimensional data in a large variety of areas including genetics, brain network analysis, social networks and climate studies. Let $G = (\mathcal{V}, \mathcal{E})$ be an undirected graph with a vertex set $\mathcal{V} = \{1, 2, \ldots, p\}$ and an edge set $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$. Let $X^0 = (X_1, X_2, ..., X_p)$ be a random vector indexed by graph's vertices. The joint distribution of $X^0$ belongs to the graphical model determined by $G$ if $X_j$ and $X_k$ are conditionally independent given all other variables whenever $j$ and $k$ are not adjacent in $G$. The graph then encodes conditional independence structure among the entries of $X^0$.

If we assume that $X^0$ is normally-distributed with a covariance matrix $\Sigma_0$, one can show that the edge structure of the graph is encoded by the precision matrix $\Theta_0 := \Sigma_0^{-1}$ (assumed to exist). If $\Theta_{ij}^0$ denotes the $(i,j)$-th entry of the matrix $\Theta_0$, it is well known that $\Theta_{ij}^0 = 0 \Leftrightarrow (i,j) \notin \mathcal{E}$. The non-zero entries in

1

the precision matrix correspond to edges in the associated graphical model and the absolute values of these entries correspond to edge weights.

Therefore to estimate the edge structure of a Gaussian graphical model, we consider the problem of estimating the precision matrix, based on a sample of $n$ independent instances $X^1, \ldots, X^n$, distributed as $X^0$. We are not only interested in point estimation, but in quantifying the uncertainty of estimation such as constructing confidence intervals and tests for edge weights. Confidence intervals and tests can be used for identifying significant variables or testing whether networks corresponding to different populations are identical.

The challenge arises due to the high-dimensional regime where the number of unknown parameters can be much larger than the number of observations $n$. It is instructive to first inspect the low-dimensional setting. In the regime when $p$ is fixed and the observations are normally distributed with $\mathbb{E}X^i = 0, i = 1, \ldots, n$, the sample covariance matrix $\hat{\Sigma} := X^T X / n$ (where $X$ is the $n \times p$ matrix of observations $X^1, \ldots, X^n$) is the maximum likelihood estimator of the covariance matrix. The inverse of the sample covariance matrix $\hat{\Theta} := \hat{\Sigma}^{-1}$ is the maximum likelihood estimator of the precision matrix. Asymptotic linearity of $\hat{\Theta}$ follows by the decomposition

$$\hat{\Theta} - \Theta_0 = -\Theta_0(\hat{\Sigma} - \Sigma_0)\Theta_0 + \text{rem}_0, \tag{1}$$

where $\text{rem}_0 := -\Theta_0(\hat{\Sigma} - \Sigma_0)(\hat{\Theta} - \Theta_0)$ is the remainder term. The term $\text{rem}_0$ can be bounded by Hölder's inequality to obtain

$$\|\text{rem}_0\|_\infty \leq \|\Theta_0(\hat{\Sigma} - \Sigma_0)\|_\infty \left\|\left\|\hat{\Theta} - \Theta_0\right\|\right\|_1,$$

where we used the notation $\|A\|_\infty = \max_{1 \leq i,j \leq p} |A_{ij}|$ for the supremum norm of a matrix $A$ and $\|A\|_1 := \max_{1 \leq j \leq p} \sum_{i=1}^p |A_{ij}|$ for the $\ell_1$-operator norm. If the fourth moments of $X^i$'s are bounded, the decomposition (1) implies rates of convergence $\|\hat{\Theta} - \Theta_0\|_\infty = \mathcal{O}_P(1/\sqrt{n})$, where $\mathcal{O}_P(1)$ denotes boundedness in probability. The remainder term then satisfies $\|\text{rem}_0\|_\infty = o_P(1/\sqrt{n})$, $o_P(1)$ denoting convergence in probability to zero. Therefore, $\hat{\Theta}$ is indeed an asymptotically linear estimator of $\Theta_0$ and in this sense, we can say it is asymptotically unbiased. Moreover, $\hat{\Theta}$ is asymptotically normal with a limiting normal distribution.

In high-dimensional settings, the sample covariance matrix does not perform well (see Johnstone (2001) and Johnstone and Lu (2009)) and if $p > n$, it is singular with probability one. Various methods have been proposed that try to reduce the variance of the sample covariance matrix at the price of introducing some bias. The idea of banding or thresholding the sample covariance matrix was studied in Bickel and Levina (2008b), Bickel and Levina (2008a) and El Karoui (2008). Methods inducing sparsity through Lasso regularization were studied by another stream of works. These can be divided into two categories: global methods and local (nodewise) methods. Global methods estimate the precision matrix typically via a regularized log-likelihood, while nodewise methods split the problem into a series of linear regressions by estimating neighbourhood of each node in the underlying graph. A popular global method is the

$\ell_1$-penalized maximum likelihood estimator, known as the graphical Lasso. Its theoretical properties were studied in a number of papers, including Yuan and Lin (2007), Friedman, Hastie, and Tibshirani (2008), Rothman, Bickel, Levina, and Zhu (2008) and Ravikumar, Raskutti, Wainwright, and Yu (2008). The local approach on estimation of precision matrices in particular includes the regression approach Meinshausen and Bühlmann (2006),Yuan (2010),Cai, Liu, and Luo (2011),Sun and Zhang (2012) which uses a Lasso-type algorithm or Dantzig selector (Candes and Tao, 2007) to estimate each column or a smaller part of the precision matrix individually.

Inference for parameters in high-dimensional undirected graphical models was studied in several papers. Multiple testing for conditional dependence in Gaussian graphical models with asymptotic control of false discovery rates was considered in Liu et al. (2013). The work Wasserman, Kolar, Rinaldo, et al. (2014) proposes methodology for inference about edge weights based on Berry-Esseen bounds and the bootstrap for certain special classes of high-dimensional graphs. Another line of work (Ren, Sun, Zhang, Zhou, et al. (2015), Janková and van de Geer (2014) and Janková and van de Geer (2016b)) proposes asymptotically normal estimators for edge weights in Gaussian graphical models based on different modifications of initial Lasso-regularized estimators. In particular, the paper Ren et al. (2015) proposes nodewise regression where each pair of variables, $(X_i, X_j)$, is regressed on all the remaining variables; this yields estimates for the parameters of the joint conditional distribution of $(X_i, X_j)$ given all the other variables. The papers Janková and van de Geer (2014) and Janková and van de Geer (2016b) propose methodology inspired by the de-biasing approach in high-dimensional linear regression that was studied in Zhang and Zhang (2014), van de Geer, Bühlmann, Ritov, and Dezeure (2014) and Javanmard and Montanari (2014). This chapter discusses and unifies the ideas from the papers Janková and van de Geer (2014) and Janková and van de Geer (2016b).

A different approach to structure learning in undirected graphical models is the Hyvärinen score matching (see e.g. Drton and Maathuis for a discussion of this approach). Methodology for asymptotically normal estimation of edge parameters in pairwise (not necessarily Gaussian) graphical models based on Hyvärinen scoring was proposed in Yu, Kolar, and Gupta (2016).

## 1.2 De-biasing regularized estimators

The idea of using regularized estimators for construction of asymptotically normal estimators is based on removing the bias associated with the penalty. Consider a real-valued loss function $\rho_\Theta$ and let $R_n(\Theta) := \sum_{i=1}^{n} \rho_\Theta(X^i)/n$ denote the average risk, given an independent sample $X^1, \ldots, X^n$. Under differentiability conditions, a regularized M-estimator $\hat{\Theta}$ based on the risk function $R_n$ can often be characterized by estimating equations

$$\dot{R}_n(\hat{\Theta}) + \xi(\hat{\Theta}) = 0, \tag{2}$$

where $\dot{R}_n$ is the gradient of $R_n$ and $\xi(\hat{\Theta})$ is a (sub-)gradient corresponding to the regularization term, evaluated at $\hat{\Theta}$. The idea is to improve on the

initial estimator by finding a root $\hat{\Theta}_{\text{de-bias}}$ closer to the solution of estimating equations without the bias term $\xi(\hat{\Theta})$, i.e. a new estimator $\hat{\Theta}_{\text{de-bias}}$ such that $\dot{R}_n(\hat{\Theta}_{\text{de-bias}}) \approx 0$. A natural way is to define a corrected estimator $\hat{\Theta}_{\text{de-bias}}$ from a linear approximation to $\dot{R}_n$

$$\dot{R}_n(\hat{\Theta}) + \ddot{R}_n(\hat{\Theta})(\hat{\Theta}_{\text{de-bias}} - \hat{\Theta}) = 0. \tag{3}$$

In high-dimensional settings, the matrix $\ddot{R}_n(\hat{\Theta})$ is typically rank deficient and thus not invertible. Suppose that we have an approximate inverse denoted by $\ddot{R}_n(\hat{\Theta})^{\text{inv}}$. Then we can approximately solve (3) for $\hat{\Theta}_{\text{de-bias}}$ to obtain

$$\hat{\Theta}_{\text{de-bias}} \approx \hat{\Theta} - \ddot{R}_n(\hat{\Theta})^{\text{inv}} \dot{R}_n(\hat{\Theta}), \tag{4}$$

provided that the remainder term is small. We refer to the step (4) as the de-biasing step since the correction term is proportional to the bias term. Generally speaking, if the initial estimator $\hat{\Theta}$ and the approximate inverse of $\ddot{R}_n(\hat{\Theta})$ are consistent in a strong-enough sense, then the new estimator $\hat{\Theta}_{\text{de-bias}}$ will be a consistent estimator of its population version $\Theta_0$ per entry at the parametric rate. The de-biasing step (4) may be viewed as one step of the Newton-Raphson scheme for numerical optimization.

In consecutive sections, we will look in detail at the bias of several particular examples of regularized estimators, including the graphical Lasso (Yuan and Lin (2007)) and nodewise Lasso (Meinshausen and Bühlmann (2006)). We now provide a unified de-biasing scheme which covers both special cases treated below (see also van de Geer (2016), Chapter 14). Suppose that a (possibly non-symmetric) estimator $\hat{\Theta}$ is available which is an approximate inverse of $\hat{\Sigma}$ in the sense that the following condition is satisfied

$$\hat{\Sigma}\hat{\Theta} - I + \eta(\hat{\Theta}) = 0, \tag{5}$$

where $\eta(\hat{\Theta})$ is a bias term. This condition in some sense corresponds to the estimating equations (2). We can express $\hat{\Theta}$ from (5) by straightforward algebra which leads to the decomposition

$$\hat{\Theta} + \hat{\Theta}^T \eta(\hat{\Theta}) - \Theta_0 = -\Theta_0(\hat{\Sigma} - \Sigma_0)\Theta_0 + \text{rem}_0 + \text{rem}_{\text{reg}}, \tag{6}$$

where

$$\text{rem}_{\text{reg}} = (\hat{\Theta} - \Theta_0)^T \eta(\hat{\Theta}).$$

Compared to the regime with $p$ fixed, there is an additional remainder $\text{rem}_{\text{reg}}$ corresponding to the bias term. Provided that the remainder terms $\text{rem}_0$ and $\text{rem}_{\text{reg}}$ are small enough, we can take as a new, de-biased estimator, $\hat{T} := \hat{\Theta} + \hat{\Theta}^T \eta(\hat{\Theta})$. The bias term $\eta(\hat{\Theta})$ can be expressed from (5) as $\eta(\hat{\Theta}) = -(\hat{\Sigma}\hat{\Theta} - I)$. Hence we obtain

$$\hat{T} = \hat{\Theta} + \hat{\Theta}^T - \hat{\Theta}^T \hat{\Sigma} \hat{\Theta}. \tag{7}$$

Bounding the remainders $\text{rem}_0$ and $\text{rem}_{\text{reg}}$ in high-dimensional settings requires more refined arguments than when $p$ is fixed. Looking at the remainder $\text{rem}_{\text{reg}}$, we can again invoke Hölder's inequality to obtain

$$\|\text{rem}_{\text{reg}}\|_\infty = \|(\hat{\Theta} - \Theta_0)^T \eta(\hat{\Theta})\|_\infty \leq \left\|\left|\hat{\Theta} - \Theta_0\right|\right\|_1 \|\eta(\hat{\Theta})\|_\infty.$$

Thus it suffices to control the rates of convergence of $\hat{\Theta}$ in $\|\|\cdot\|\|_1$-norm and control the absolute size of entries of the bias term.

Provided that the remainders are of small order $1/\sqrt{n}$ in probability, asymptotic normality per elements of $\hat{T}$ is a consequence of asymptotic linearity and can be established under tail conditions on $X^i$'s, by applying the Lindeberg's central limit theorem.

## 1.3 Graphical Lasso

If the observations are independent $\mathcal{N}(0, \Sigma_0)$-distributed, the log-likelihood function is proportional to

$$\ell(\Theta) := \text{tr}(\hat{\Sigma}\Theta) - \log \det(\Theta).$$

The graphical Lasso (see Yuan and Lin (2007), d'Aspremont, Banerjee, and El Ghaoui (2008), Friedman et al. (2008)) is based on the Gaussian log-likelihood function but regularizes it via an $\ell_1$-norm penalty on the off-diagonal elements of the precision matrix. The diagonal elements of the precision matrix correspond to certain partial variances and thus should not be penalized. The graphical Lasso is defined by

$$\hat{\Theta} = \text{argmin}_{\Theta=\Theta^T, \Theta\succ 0}\text{tr}(\hat{\Sigma}\Theta) - \log \det(\Theta) + \lambda\|\Theta^-\|_1, \tag{8}$$

where $\lambda$ is non-negative tuning parameter and we optimize over symmetric positive definite matrices, denoted by $\succ$. Here $\Theta^-$ represents the matrix obtained by setting the diagonal elements of $\Theta$ to zero and $\|\Theta^-\|_1$ is the $\ell_1$-norm of the vectorized version of $\Theta^-$. The usefulness of the graphical Lasso is not limited only to Gaussian settings; the theoretical results below show that it performs well as an estimator of the precision matrix in a large class of non-Gaussian settings.

A disadvantage of the graphical Lasso (8) is that the penalization does not take into account that the variables have in general a different scaling. To take these differences in the variances into account, we define a modified graphical Lasso with a weighted penalty. To this end, let $\hat{W}^2 := \text{diag}(\hat{\Sigma})$ be the diagonal matrix obtained from the diagonal of $\hat{\Sigma}$. We let

$$\hat{\Theta}_{\text{w}} = \text{argmin}_{\Theta=\Theta^T, \Theta\succ 0}\text{tr}(\hat{\Sigma}\Theta) - \log \det(\Theta) + \sum_{i\neq j} \hat{W}_{ii}\hat{W}_{jj}|\Theta_{ij}|. \tag{9}$$

The weighted graphical Lasso $\hat{\Theta}_{\text{w}}$ is related to a graphical Lasso based on the sample correlation matrix $\hat{R} := \hat{W}^{-1}\hat{\Sigma}\hat{W}^{-1}$. To clarify the connection, we define

$$\hat{\Theta}_{\text{norm}} = \text{argmin}_{\Theta=\Theta^T, \Theta\succ 0}\text{tr}(\hat{R}\Theta) - \log \det(\Theta) + \|\Theta^-\|_1. \tag{10}$$

Then it holds that $\hat{\Theta}_{\text{w}} = \hat{W}^{-1}\hat{\Theta}_{\text{norm}}\hat{W}^{-1}$. The estimator $\hat{\Theta}_{\text{norm}}$ is of independent interest, if the parameter of interest is the inverse correlation matrix rather than the precision matrix. The estimators $\hat{\Theta}_{\text{w}}$ and $\hat{\Theta}_{\text{norm}}$ are also useful from a theoretical perspective as will be shown in the sequel.

We now apply the de-biasing ideas of Section 1.2 to the graphical Lasso estimators defined above, demonstrating the procedure on $\hat{\Theta}$. By definition, the graphical Lasso is invertible, and the Karush-Kuhn-Tucker (KKT) conditions yield

$$\hat{\Sigma} - \hat{\Theta}^{-1} + \lambda\hat{Z} = 0,$$

where

$$\hat{Z}_{ij} = \text{sign}(\hat{\Theta}_{ij}) \ \text{ if } \hat{\Theta}_{ij} \neq 0, \quad \text{ and } \quad \|\hat{Z}\|_\infty \leq 1.$$

Multiplying by $\hat{\Theta}$, we obtain

$$\hat{\Sigma}\hat{\Theta} - I + \lambda\hat{Z}\hat{\Theta} = 0.$$

In line with Section 1.2 above, this implies the decomposition

$$\hat{\Theta} + \hat{\Theta}^T\eta(\hat{\Theta}) - \Theta_0 = -\Theta_0(\hat{\Sigma} - \Sigma_0)\Theta_0 + \text{rem}_0 + \text{rem}_{\text{reg}},$$

with the bias term $\eta(\hat{\Theta}) = \lambda\hat{Z}\hat{\Theta}$. To control the remainder terms $\text{rem}_0$ and $\text{rem}_{\text{reg}}$, we need bounds for the $\ell_1$-error of $\hat{\Theta}$ and to control the bias term, it is sufficient to control the upper bound $\|\eta(\hat{\Theta})\|_\infty = \|\lambda\hat{Z}\hat{\Theta}\|_\infty \leq \lambda\left|\!\left|\!\left|\hat{\Theta}\right|\!\right|\!\right|_1$.

### Oracle bounds

Oracle results for the graphical Lasso were studied in Rothman et al. (2008) under sparsity conditions and mild regularity conditions. In Ravikumar et al. (2008), stronger results were derived under stronger regularity conditions (and weaker sparsity conditions). Here we revisit these results and provide several extensions.

We summarize the main theoretical conditions which require boundedness of the eigenvalues of the true precision matrix and certain tail conditions.

**Condition A1** (Bounded spectrum). *The precision matrix $\Theta_0 := \Sigma_0^{-1}$ exists and there exists a universal constant $L \geq 1$ such that*

$$1/L \leq \Lambda_{\min}(\Theta_0) \leq \Lambda_{\max}(\Theta_0) \leq L.$$

**Condition A2** (Sub-Gaussianity). *The observations $X^i, i = 1, \ldots, n$, are uniformly sub-Gaussian vectors, i.e. there exists a universal constant $K > 0$ such that for every $\alpha \in \mathbb{R}^p$, $\|\alpha\|_2 = 1$ it holds*

$$\mathbb{E}\exp\left(|\alpha^T X^i|^2/K^2\right) \leq 2 \quad (i = 1, \ldots, n).$$

Under Condition A2, the Bernstein inequality implies concentration results for $\hat{\Sigma}$ as formulated in Lemma 1 below. The proof is omitted and may be found in Bühlmann and van de Geer (2011) (Lemma 14.13). We denote the Euclidean norm by $\|\cdot\|_2$ and the $i$-th column of a matrix $A$ by $A_i$.

**Lemma 1.** *Assume Condition A2 and that for non-random matrices $A, B \in \mathbb{R}^{p \times p}$ it holds that $\|A_i\|_2 \leq M$ and $\|B_i\|_2 \leq M$ for all $i = 1, \ldots, p$. Then for all $t > 0$, with probability at least $1 - e^{-nt}$ it holds that*

$$\|A^T(\hat{\Sigma} - \Sigma_0)B\|_\infty / (2M^2K^2) \leq t + \sqrt{2t} + \sqrt{\frac{2\log(2p^2)}{n}} + \frac{\log(2p^2)}{n}.$$

To derive oracle bounds for the graphical Lasso, we rely on certain sparsity conditions on the entries of the true precision matrix. To this end, we define for $j = 1, \ldots, p$,

$$D_j := \{(i,j) : \Theta_{ij}^0 \neq 0, i \neq j\}, \quad d_j := \text{card}(D_j), \quad d := \max_{j=1,\ldots,p} |d_j|.$$

The quantity $d_j$ is then the degree of a node $X_j$ and $d$ corresponds to the maximum vertex degree in the associated graphical model (excluding vertex self-loops). Furthermore, we define

$$S := \bigcup_{j=1}^{p} D_j, \quad s := \sum_{j=1}^{p} d_j,$$

thus $S$ denotes the overall off-diagonal sparsity pattern and $s$ is the overall number of edges (excluding self-loops).

The following theorem is an extension of the result for the graphical Lasso in Rothman et al. (2008) to the $\ell_1$-norm. The paper Rothman et al. (2008) derives rates in Frobenius norm $\|\cdot\|_F$, which is defined as $\|A\|_F^2 := \sum_{i,j} |A_{ij}^2|$ for a matrix $A$.

**Theorem 1** (Regime $p \ll n$). *Let $\hat{\Theta}$ be the minimizer defined by* (8). *Assume Conditions A1 and A2. Then for $\lambda$ satisfying $2\lambda_0 \leq \lambda \leq 1/(8Lc_L)$ and $8c_L^2 s\lambda^2 + 8c_L p\lambda_0^2 \leq \lambda_0/(2L)$, on the set $\|\hat{\Sigma} - \Sigma_0\|_\infty \leq \lambda_0$, it holds that*

$$\|\hat{\Theta} - \Theta_0\|_F^2/c_L + \lambda\|\hat{\Theta}^- - \Theta_0^-\|_1 \leq 8c_L^2 s\lambda^2 + 8c_L p\lambda_0^2,$$

*and*

$$\left\|\left|\hat{\Theta} - \Theta_0\right|\right\|_1 \leq 16c_L^2(p+s)\lambda,$$

*where $c_L = 8L^2$.*

The slow rate in the result above arises from the part of the estimation error $\text{tr}[(\hat{\Sigma} - \Sigma_0)(\hat{\Theta} - \Theta_0)]$ which is related to the diagonal elements of the precision matrix. However, proper normalizing removes this part of the estimation error.

The following theorem derives an extension of Rothman et al. (2008) for the normalized graphical Lasso $\hat{\Theta}_{\text{norm}}$. Denote the true correlation matrix by $R_0$ and let $K_0 := R_0^{-1}$ denote the inverse correlation matrix.

**Theorem 2** (Regime $p \gg n$). *Assume that Conditions A1 and A2 hold. Then for $\lambda$ satisfying $2\lambda_0 \leq \lambda \leq 1/(8L^2)$ and $8c_L^2 s\lambda^2 \leq \lambda_0/(2L)$, on the set $\|\hat{R} - R_0\|_\infty \leq \lambda_0$ it holds for some constant $C_L > 0$ that*

$$\|\hat{\Theta}_{\text{norm}} - K_0\|_F^2 + \lambda\|\hat{\Theta}_{\text{norm}}^- - K_0^-\|_1 \leq 8c_L^2 s\lambda^2,$$

$$\left\| \hat{\Theta}_{\mathrm{norm}} - K_0 \right\|_1 \leq 8 c_L s \lambda^2 + 8 c_L^2 s \lambda.$$

$$\left\| \hat{\Theta}_{\mathrm{w}} - \Theta_0 \right\|_1 \leq C_L s \lambda,$$

where $c_L = 8 L^2$.

Using the normalized graphical Lasso leads to faster rates in Frobenius norm and $\ell_1$-norm as shown above. The rates for $\hat{\Theta}_{\mathrm{w}}$ in $\|\cdot\|_1$-norm can be then established immediately. To derive a high-probability bound for $\|\hat{R} - R_0\|_\infty$, we may apply Lemma 1 together with Hölder's inequality to obtain $\|\hat{R} - R_0\|_\infty = \mathcal{O}_P(\sqrt{\log p / n})$. Hence, $\left\| \hat{\Theta}_{\mathrm{w}} - \Theta_0 \right\|_1 = \mathcal{O}_P(s \sqrt{\log p / n})$.

**Remark 1.** The above result requires a strong condition on the sparsity in $\Theta_0$, i.e. there can only be very few non-zero coefficients due to the restriction $8 c_L^2 s \lambda^2 \leq \lambda_0 / (2L)$. This condition guarantees that a margin condition is satisfied. An example of a graph satisfying the condition is a star graph with order $\sqrt{n}$ edges.

### Asymptotic normality

Once oracle results in $\ell_1$-norm are available, we can easily obtain results on asymptotic normality of the de-biased estimator $2\hat{\Theta} - \hat{\Theta}\hat{\Sigma}\hat{\Theta}$ for the graphical Lasso and its weighted version. We denote the asymptotic variance of the de-biased estimator by $\sigma_{ij}^2 := n\mathrm{var}((\Theta_i^0)^T \hat{\Sigma} \Theta_j^0)$. The arrow $\rightsquigarrow$ denotes convergence in distribution and for a matrix $A$ we denote $(A)_{ij}$ its $(i,j)$-entry.

**Theorem 3** (Regime $p \ll n$). *Assume Conditions A1, A2, $\lambda \asymp \sqrt{\log p / n}$ and that $(p + s)\sqrt{d} = o(\sqrt{n} / \log p)$. Then it holds that*

$$2\hat{\Theta} - \hat{\Theta}\hat{\Sigma}\hat{\Theta} - \Theta_0 = -\Theta_0(\hat{\Sigma} - \Sigma_0)\Theta_0 + \mathrm{rem}, \tag{11}$$

*where*

$$\|\mathrm{rem}\|_\infty = \mathcal{O}_P\left(24(8L^2)^2 L(p + s)\sqrt{d+1}\lambda^2\right) = o_P(1/\sqrt{n}).$$

*Moreover, for $i, j = 1, \ldots, p$,*

$$\sqrt{n}(2\hat{\Theta} - \hat{\Theta}\hat{\Sigma}\hat{\Theta} - \Theta_0)_{ij}/\sigma_{ij} \rightsquigarrow \mathcal{N}(0, 1).$$

The result of Theorem 3 gives us tools to construct approximate confidence intervals and tests for individual entries of $\Theta_0$. However, we need a consistent estimator of the asymptotic variance $\sigma_{ij}$. For the Gaussian case, one may take $\hat{\sigma}_{ij}^2 := \hat{\Theta}_{ii}\hat{\Theta}_{jj} + \hat{\Theta}_{ij}^2$. We omit the proof of consistency of $\hat{\sigma}_{ij}$ and point the reader to Janková and van de Geer (2016b), where other distributions are treated as well. Moreover, Theorem 3 implies parametric rates of convergence for estimation of individual entries and a rate of order $\sqrt{\log p / n}$ for the error in supremum norm. Theorem 3 requires a stronger sparsity condition that the corresponding oracle-type inequality in Theorem 1. This is to be expected as will be argued in Section 1.7.

Using the weighted graphical Lasso, the results of Theorem 3 can be established under weaker conditions as shown in the following theorem.

**Theorem 4** (Regime $p \gg n$). *Assume Conditions A1, A2 and $s\sqrt{d} = o(\sqrt{n}/\log p)$. Then for $\lambda \asymp \sqrt{\log p/n}$, the asymptotic linearity (11) holds with $\hat{\Theta}_{\mathrm{w}}$, where*

$$\|\text{rem}\|_\infty = \mathcal{O}_P\left(12(8L^2)^2 s\sqrt{d+1}\lambda^2\right) = o_P(1/\sqrt{n}).$$

*Moreover, for $i, j = 1, \ldots, p$, $\sqrt{n}(2\hat{\Theta}_{\mathrm{w}} + \hat{\Theta}_{\mathrm{w}}\hat{\Sigma}\hat{\Theta}_{\mathrm{w}} - \Theta_0)_{ij}/\sigma_{ij} \rightsquigarrow \mathcal{N}(0,1)$.*

If the parameter of interest is instead the inverse correlation matrix, we formulate a partial result below.

**Proposition 1** (Regime $p \gg n$). *Assume Conditions A1, A2, $\lambda \asymp \sqrt{\log p/n}$ and that $s\sqrt{d} = o(\sqrt{n}/\log p)$. Then it holds*

$$2\hat{\Theta}_{\mathrm{norm}} - \hat{\Theta}_{\mathrm{norm}}\hat{R}\hat{\Theta}_{\mathrm{norm}} - K_0 = -K_0(\hat{R} - R_0)K_0 + \text{rem},$$

$$\|\text{rem}\|_\infty = \mathcal{O}_P\left(12(8L^2)^2 Ls\sqrt{d+1}\lambda^2\right) = o_P(1/\sqrt{n}).$$

To claim asymptotic normality of $K_0(\hat{R} - R_0)K_0$ per entry would require extensions of central limit theorems to high-dimensional settings (see Chernozhukov, Chetverikov, and Kato (2014)) and an extension of the $\delta$-method. We do not study these extensions in the present work. To give a glimpse, in the regime when $p$ is fixed, by the central limit theorem it follows that $\sqrt{n}(\hat{\Sigma} - \Sigma_0) \rightsquigarrow \mathcal{N}_{p^2}(0, C)$, where $C$ is the asymptotic covariance matrix. Then we may apply the $\delta$-method with $h_{ij}(\Sigma) := (K_i^0)^T \text{diag}(\Sigma)^{-1/2} \Sigma \text{diag}(\Sigma)^{-1/2} K_j^0$ to obtain $\sqrt{n}(h_{ij}(\hat{\Sigma}) - h_{ij}(\Sigma_0)) \rightsquigarrow \mathcal{N}(0, \dot{h}(\Sigma_0)^T C \dot{h}(\Sigma_0))$.

Finally, we show that the sparsity conditions in the above results may be further relaxed under a stronger regularity condition on the true precision matrix. We provide here a simplified version of the result in Janková and van de Geer (2014) which assumes an irrepresentability condition on the true precision matrix. Let $\kappa_{\Sigma_0}$ be the $\ell_\infty$-operator norm of the true covariance matrix $\Sigma_0$, i.e. $\kappa_{\Sigma_0} = \|\Sigma_0\|_1$. Let $H_0$ be the Hessian of the expected Gaussian log-likelihood evaluated at $\Theta_0$, i.e. $H_0 = \Sigma_0 \otimes \Sigma_0$. When calculating the Hessian matrix, we treat the precision matrix as non-symmetric; this will allow us to accommodate non-symmetric estimators as well. For any two subsets $T$ and $T'$ of $\mathcal{V} \times \mathcal{V}$, we use $H_{TT'}^0$ to denote the $|T| \times |T'|$ matrix with rows and columns of $H_0$ indexed by $T$ and $T'$ respectively. Define $\kappa_{H_0} = \left\|(H_{SS}^0)^{-1}\right\|_1$.

**Condition A3.** *(Irrepresentability condition) There exists $\alpha \in (0, 1]$ such that $\max_{e \in S^c} \|H_{eS}^0 (H_{SS}^0)^{-1}\|_1 \leq 1 - \alpha$, where $S^c$ is the complement of $S$.*

Condition A3 is an analogy of the irrepresentable condition for variable selection in linear regression (see van de Geer and Bühlmann (2009)). If we define the zero-mean edge random variables as $Y_{(i,j)} := X_i X_j - \mathbb{E}(X_i X_j)$, then the matrix $H_0$ corresponds to covariances of the edge variables, in particular $H_{(i,j),(k,l)}^0 + H_{(j,i),(k,l)}^0 = \text{cov}(Y_{(i,j)}, Y_{(k,l)})$. Condition A3 means that we require that no edge variable $Y_{(j,k)}$, which is not included in the edge set $S$, is highly correlated with variables in the edge set (see Ravikumar et al. (2008)). The

parameter $\alpha$ then is a measure of this correlation with the correlation growing when $\alpha \to 0$. Some examples of matrices satisfying the irrepresentable condition may be found in Janková and van de Geer (2016b).

**Theorem 5** (Regime $p \gg n$). *Assume that Conditions A1, A2 and A3 are satisfied, $d = o(\sqrt{n}/\log p)$, $\kappa_{\Sigma_0} = \mathcal{O}(1)$ and $\kappa_{H_0} = \mathcal{O}(1)$. Then for $\lambda \asymp \sqrt{\log p/n}$, the asymptotic linearity (11) holds with $\hat{\Theta}$, where $\|\mathrm{rem}\|_\infty = \mathcal{O}_P(d \log p/n) = o_P(1/\sqrt{n})$. Moreover,*

$$\sqrt{n}(2\hat{\Theta} - \hat{\Theta}\hat{\Sigma}\hat{\Theta} - \Theta_0)_{ij}/\sigma_{ij} \rightsquigarrow \mathcal{N}(0,1).$$

The proof of Theorem 5 may be found in Janková and van de Geer (2014). We remark that under the irrepresentability condition, one can show that $|\hat{\Theta}_{ij} - \Theta_{ij}^0| = \mathcal{O}_P(1/\sqrt{n})$ (see Ravikumar et al. (2008)). This means that one could use the graphical Lasso itself to construct confidence intervals of asymptotically optimal (parametric) size. However, this holds under the strong irrepresentability condition which is often violated in practice.

Comparing the results obtained for the (weighted) graphical Lasso, we see that the strongest result was attained under the irrepresentable condition and under the sparsity condition $d = o(\sqrt{n}/\log p)$. An analogous result has not yet been obtained for the graphical Lasso without the irrepresentable condition (under the same sparsity condition). However, without the irrepresentable condition, we showed the same results for the weighted graphical Lasso under the sparsity condition $s\sqrt{d} = o(\sqrt{n}/\log p)$. In the next section, we will consider a procedure based on pseudo-likelihood, for which we can derive identical results under weaker conditions, namely under the sparsity condition $d = o(\sqrt{n}/\log p)$ and under the Conditions A1 and A2.

## 1.4 Nodewise square-root Lasso

An alternative approach to estimate the precision matrix is based on linear regression. The idea of nodewise Lasso is to estimate each column of the precision matrix by doing a projection of every column of the design matrix on all the remaining columns. While this is a pseudo-likelihood method, the decoupling into linear regressions gains more flexibility in estimating the individual scaling levels compared to the graphical Lasso which aims to estimate all the parameters simultaneously. Moreover, by splitting the problem up into a series of linear regressions, the computational burden is reduced compared to the graphical Lasso.

In low-dimensional settings, regressing each column of the design matrix on all the other columns would simply recover the inverse of the sample covariance matrix $\hat{\Sigma}$. However, due to the high-dimensionality of our setting, the matrix $\hat{\Sigma}$ is not invertible and we can only do approximate projections. If we assume sparsity in the precision matrix (and thus also in the partial correlations), this idea can be effectively carried out using the square-root Lasso (Belloni, Chernozhukov, and Wang (2011)).

The theoretical motivation can be understood in greater detail from the population version of the method. For each $j = 1, \ldots, p$, we define the vector of partial correlations $\gamma_j^0 = \{\gamma_{j,k}^0, k \neq j\}$ as follows

$$\gamma_j^0 := \mathrm{argmin}_{\gamma \in \mathbb{R}^{p-1}} \mathbb{E} \|X_j - X_{-j}\gamma\|_2^2 / n, \tag{12}$$

and we denote the noise level by $\tau_j^2 = \mathbb{E} \|X_j - X_{-j}\gamma_j^0\|_2^2 / n$. Then one may show that the $j$-th column of $\Theta_0$ can be recovered from the partial correlations $\gamma_j^0$ and the noise level $\tau_j$ using the following identity: $\Theta_j^0 = (-\gamma_{j,1}, \ldots, -\gamma_{j,j-1}, 1, -\gamma_{j,j+1}, \ldots, -\gamma_{j,p})^T / \tau_j^2$.

Hence, the idea is to define for each $j = 1, \ldots, p$ the estimators of the partial correlations, $\hat{\gamma}_j = \{\hat{\gamma}_{j,k}, k = 1, \ldots, p, j \neq k\} \in \mathbb{R}^{p-1}$ using, for instance, the square-root Lasso with weighted penalty,

$$\hat{\gamma}_j := \mathrm{argmin}_{\gamma \in \mathbb{R}^{p-1}} \|X_j - X_{-j}\gamma\|_2 / n + 2\lambda \|\hat{W}_{-j}\gamma\|_1, \tag{13}$$

where by $A_{-j}$ we denote the matrix $A$ without its $j$-th column. We further define estimators of the noise level

$$\hat{\tau}_j^2 := \|X_j - X_{-j}\hat{\gamma}_j\|_2^2 / n, \quad \tilde{\tau}_j^2 := \hat{\tau}_j^2 + \lambda \hat{\tau}_j \|\hat{\gamma}_j\|_1,$$

for $j = 1, \ldots, p$. Finally, we define the nodewise square-root Lasso estimator

$$\hat{\Theta} := \begin{pmatrix} 1/\tilde{\tau}_1^2 & -\tilde{\gamma}_{1,2}/\tilde{\tau}_1^2 & \cdots & -\tilde{\gamma}_{1,p}/\tilde{\tau}_1^2 \\ -\tilde{\gamma}_{2,1}/\tilde{\tau}_2^2 & 1/\tilde{\tau}_2^2 & \cdots & -\tilde{\gamma}_{2,p}/\tilde{\tau}_2^2 \\ \vdots & \vdots & \ddots & \vdots \\ -\tilde{\gamma}_{p,1}/\tilde{\tau}_p^2 & \cdots & -\tilde{\gamma}_{p,p-1}/\tilde{\tau}_p^2 & 1/\tilde{\tau}_p^2 \end{pmatrix} \tag{14}$$

An equivalent way of formulating the definitions above is

$$(\hat{\gamma}_j, \hat{\tau}_j) = \mathrm{argmin}_{\gamma \in \mathbb{R}^{p-1}, \tau \in \mathbb{R}} \|X_j - X_{-j}\gamma\|_2^2 / n / (2\tau) + \tau/2 + 2\lambda \|\hat{W}_{-j}\gamma\|_1. \tag{15}$$

Alternative versions of the above estimator were considered in the literature. One may use the Lasso (Tibshirani (1996)) instead of the square-root Lasso (as in Janková and van de Geer (2016b)) or the Dantzig selector (see van de Geer (2016)). Furthermore, one may define the nodewise square-root Lasso with $\hat{\tau}_j$ in place of $\tilde{\tau}_j$.

The properties of the column estimator $\hat{\Theta}_j$ were studied in several papers (following Meinshausen and Bühlmann (2006)) and it has been shown to enjoy oracle properties under the Conditions A1, A2 and under the sparsity condition $d = \mathcal{O}(n/\log p)$ (see van de Geer et al. (2014), where a similar version was considered).

In line with Section 1.2, we consider a de-biased version of the nodewise square-root Lasso estimator. The KKT conditions for the optimization problem (13) give

$$-\hat{\tau}_j X_{-j}^T (X_j - X_{-j}\hat{\gamma}_j) / n + \lambda \hat{\kappa}_j = 0, \tag{16}$$

11

for $j = 1, \ldots, p$, where $\hat{\kappa}_j$ is an element of the sub-differential of the function $\gamma_j \mapsto \|\gamma_j\|_1$ at $\hat{\gamma}_j$, i.e. for $k \in \{1, \ldots, p\} \setminus \{j\}$,

$$\hat{\kappa}_{j,k} = \text{sign}(\hat{\gamma}_{j,k}) \text{ if } \hat{\gamma}_{j,k} \neq 0, \quad \text{and} \quad \|\hat{\kappa}_j\|_\infty \leq 1.$$

If we define $\hat{Z}_j$ to be a $p \times 1$ vector

$$\hat{Z}_j := (\hat{\kappa}_{j,1}, \ldots, \hat{\kappa}_{j,j-1}, 0, \hat{\kappa}_{j,j+1}, \ldots, \hat{\kappa}_{j,p}),$$

then the KKT conditions may be equivalently stated as follows

$$\hat{\Sigma}\hat{\Theta}_j - e_j - \lambda \frac{\hat{\tau}_j}{\tilde{\tau}_j^2} \hat{Z}_j = 0.$$

Let $\hat{Z}$ be a matrix with columns $\hat{Z}_j$ for $j = 1, \ldots, p$, $\hat{\tau}$ be a diagonal matrix with elements $(\hat{\tau}_1, \ldots, \hat{\tau}_p)$ and similarly $\tilde{\tau} := \text{diag}(\tilde{\tau}_1, \ldots, \tilde{\tau}_p)$. As in Section 1.2, this yields the decomposition (6) with a bias term $\eta(\hat{\Theta}) := \hat{Z}\Lambda\hat{\tau}\tilde{\tau}^{-2}$. The bias term can then be controlled as

$$\|\eta(\hat{\Theta})\|_\infty \leq \lambda\|\hat{\tau}\|_\infty\|\tilde{\tau}^{-2}\|_\infty\|\hat{Z}\|_\infty \leq \lambda \max_{1 \leq j \leq p} \hat{\tau}_j/\tilde{\tau}_j^2.$$

**Theorem 6** (Regime $p \gg n$). *Suppose that Conditions A1, A2 are satisfied and $d = o(\sqrt{n}/\log p)$. Let $\hat{\Theta}_{\text{node}}$ be the estimator defined in (14) and let $\lambda \asymp \sqrt{\log p/n}$. Then it holds*

$$\hat{\Theta} + \hat{\Theta}^T - \hat{\Theta}^T\hat{\Sigma}\hat{\Theta} - \Theta_0 = -\Theta_0(\hat{\Sigma} - \Sigma_0)\Theta_0 + \text{rem},$$

*where $\|\text{rem}\|_\infty = \mathcal{O}_P(d\lambda^2) = o_P(1/\sqrt{n})$. Moreover,*

$$\sqrt{n}(\hat{\Theta} + \hat{\Theta}^T - \hat{\Theta}^T\hat{\Sigma}\hat{\Theta} - \Theta_0)_{ij}/\sigma_{ij} \rightsquigarrow \mathcal{N}(0, 1).$$

When the parameter of interest is the inverse correlation matrix, we can use the normalized version of the nodewise square-root Lasso and we obtain an analogous result.

**Proposition 2** (Regime $p \gg n$). *Suppose that Conditions A1, A2 are satisfied, let $\lambda \asymp \sqrt{\log p/n}$ and $d = o(\sqrt{n}/\log p)$. Then*

$$\hat{\Theta}_{\text{norm}} + \hat{\Theta}_{\text{norm}}^T - \hat{\Theta}_{\text{norm}}^T\hat{R}\hat{\Theta}_{\text{norm}} - \Theta_0 = -K_0(\hat{R} - R_0)K_0 + \text{rem},$$

*where $\|\text{rem}\|_\infty = o_P(1/\sqrt{n})$.*

## 1.5 Computational view

For the nodewise square-root Lasso, we need to solve $p$ square-root Lasso regressions, which can be efficiently handled using interior-point methods with polynomial computational time or first-order methods (see Belloni et al. (2011)). Alternatively to nodewise square-root Lasso, the nodewise Lasso studied in Janková

12

and van de Geer (2016b) may be used, which requires selection of a tuning parameter for each of the $p$ regressions. This can be achieved e.g. by cross-validation and can be implemented efficiently using parallel methods (Efron, Hastie, Johnstone, and Tibshirani (2004)). The graphical Lasso presents a more computationally challenging problem; we refer the reader to e.g. Mazumder and Hastie (2012). The computation of the de-biased estimator itself only involves simple matrix addition and multiplication.

## 1.6   Simulation results

We consider a setting with $n$ independent observations generated from $\mathcal{N}_p(0, \Theta_0^{-1})$, where the precision matrix $\Theta_0$ follows one of the three models:

1. Model 1: $\Theta_0$ has two blocks of equal size and each block is a five-diagonal matrix with elements $(1, 0.5, 0.4)$ and $(2, 1, 0.6)$, respectively.

2. Model 2:  $\Theta_0$ is a sparse precision matrix generated using the R package GGMselect (using the function simulateGraph() with parameter 0.07). The matrix was converted to a correlation matrix with the function cov2cor().

3. Model 3: $\Theta_{ij}^0 = 0.5^{|i-j|}$, $i, j = 1, \ldots, p$.

We consider 6 different methods: the de-biased estimator based on the

(1) graphical Lasso (glasso)

(2) weighted graphical Lasso (glasso-weigh),

(3) nodewise square-root Lasso (node-sqrt) as defined above,

(4) nodewise square-root Lasso with alternative $\tilde{\tau}$ as in Sun and Zhang (2012) (node-sqrt-tau)

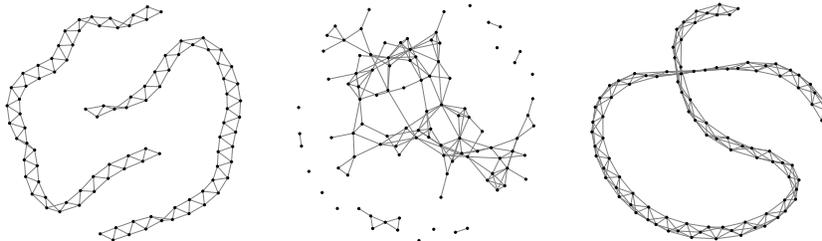(5) nodewise Lasso as in Janková and van de Geer (2016b) (node)

and we also consider

(6) the maximum likelihood estimator (MLE).

Furthermore, as a benchmark we report the oracle estimator (oracle) which applies maximum likelihood using the knowledge of true zeros in the precision matrix. We also report the target coverage and the efficient asymptotic variance of asymptotically regular estimators (see Janková and van de Geer (2016a)) (perfect).

For the graphical Lasso (1) and the weighted graphical Lasso (2) we choose the tuning parameter by maximizing the likelihood on a validation data set (a new data of the size $n$). For methods (3), (4) and (5), the universal choice $\sqrt{\log p/n}$ is used.

We display results on confidence intervals for nominal coverage 95% and for normally distributed observations in Tables 1, ??, 2 and 3. For other nominal coverages, we obtain similar performance (these results are not reported). For other than Gaussian distributions, we refer the reader to the simulation results in

Figure 1: Visualization of graphical models used in simulations. Models 1,2 and 3 from left to right. For Model 3, we only plot edges with a weight greater than 0.1.



Janková and van de Geer (2016b). Firstly, the results of the simulations suggest that the de-biased estimators perform significantly better than the maximum likelihood estimator even though $p < n$ and secondly, the nodewise methods seem to outperform the graphical lasso methods in our settings.

## 1.7 Discussion

We have shown several constructions of asymptotically linear estimators of the precision matrix (and inverse correlation matrix) based on regularized estimators, which immediately lead to inference in Gaussian graphical models. Efficient algorithms are available for both methods as discussed in Section 1.5. The constructed estimators achieve entrywise estimation at the parametric rate and a rate of convergence of order $\sqrt{\log p/n}$ in supremum norm.

To provide a brief comparison of the two methods analyzed above, both theoretical and computational results seem in favor of the de-sparsified nodewise Lasso. Theoretical results for nodewise Lasso in the regime $p \gg n$ only need the mild conditions Conditions A1, A2 and $d = o(\sqrt{n}/\log p)$ and are uniform over the considered model. Moreover, the de-sparsified nodewise Lasso may be thresholded again to yield recovery of the set $S$ with no false positives, and under a beta-min type condition, exact recovery of the set $S$, asymptotically, with high probability. The graphical Lasso requires that we impose the strong irrepresentability condition in the high-dimensional regime. However, the graphical Lasso might be preferred on the grounds that it does not decouple the likelihood. Moreover, the graphical Lasso estimator is always strictly positive definite and thus yields an estimator of the covariance matrix as well. The invertibility of the nodewise Lasso has not yet been explored.

We remark that the sparsity condition $d = o(\sqrt{n}/\log p)$ implied by our analysis is stronger than the condition needed for oracle inequalities and recovery, namely $d = o(n/\log p)$. However, one can show that this sparsity condition is essentially necessary for asymptotically normal estimation. This follows by inspection of the minimax rates (see Ren et al. (2015)).

Table 1: Average coverages and lengths of confidence intervals over the active set $S_0 := S \cup \{1, \ldots, p\}$ and its complement $S_0^c$, over 100 realizations. The average value of the tuning parameters is reported in the last column. The benchmark "estimators" are labeled by a star. The significance level is 0.05.

Model 1: Block 1: $(1, 0.5, 0.4)$, Block 2: $(2, 1, 0.6)$

$p = 100, n = 200$

|   | Method | Coverage $S_0$ | Coverage $S_0^c$ | Length $S_0$ | Length $S_0^c$ | Average $\lambda$ |
|---|--------|------|------|------|------|-----------|
| 1 | glasso | 77.19 | 98.07 | 0.36 | 0.32 | 0.088 |
| 2 | glasso-weigh | 35.02 | 98.65 | 0.31 | 0.27 | 0.088 |
| 3 | node-sqrt | 89.92 | 94.02 | 0.48 | 0.42 | 0.152 |
| 4 | node-sqrt-tau | 83.48 | 97.40 | 0.38 | 0.33 | 0.152 |
| 5 | node | 90.58 | 96.77 | 0.41 | 0.35 | 0.152 |
| 6 | MLE | 20.92 | 84.27 | 0.97 | 0.81 | - |
| 7 | oracle* | 94.95 | - | 0.49 | 0.40 | - |
| 8 | perfect* | 95.00 | 95.00 | 0.48 | 0.40 | - |

$p = 100, n = 400$

|   | Method | Coverage $S_0$ | Coverage $S_0^c$ | Length $S_0$ | Length $S_0^c$ | Average $\lambda$ |
|---|--------|------|------|------|------|-----------|
| 1 | glasso | 84.28 | 97.53 | 0.27 | 0.23 | 0.049 |
| 2 | glasso-weigh | 46.22 | 98.41 | 0.24 | 0.20 | 0.049 |
| 3 | node-sqrt | 91.57 | 94.40 | 0.34 | 0.29 | 0.107 |
| 4 | node-sqrt-tau | 87.11 | 97.13 | 0.28 | 0.24 | 0.107 |
| 5 | node | 91.48 | 96.40 | 0.30 | 0.25 | 0.107 |
| 6 | MLE | 41.42 | 91.29 | 0.41 | 0.38 | - |
| 7 | oracle* | 94.87 | - | 0.34 | 0.29 | - |
| 8 | perfect* | 95.00 | 95.00 | 0.34 | 0.29 | - |

15

Table 2: Average coverages and lengths of confidence intervals over the active set $S_0 := S \cup \{1, \ldots, p\}$ and its complement $S_0^c$, over 100 realizations. The average value of the tuning parameters is reported in the last column. The benchmark "estimators" are labeled by a star. The significance level is 0.05.

Model 2

$p = 100, n = 400$

|   | Method | Coverage | | Length | | Average $\lambda$ |
|---|--------|---------|---------|--------|--------|-----------|
|   |        | $S_0$   | $S_0^c$ | $S_0$  | $S_0^c$ |           |
| 1 | glasso | 64.17 | 98.65 | 0.16 | 0.15 | 0.067 |
| 2 | glasso-weigh | 16.80 | 98.56 | 0.05 | 0.05 | 0.067 |
| 3 | node-sqrt | 87.23 | 94.43 | 0.24 | 0.21 | 0.107 |
| 4 | node-sqrt-tau | 89.81 | 97.23 | 0.20 | 0.18 | 0.107 |
| 5 | node | 38.19 | 99.07 | 0.10 | 0.10 | 0.107 |
| 6 | MLE | 50.98 | 91.22 | 0.30 | 0.26 | - |
| 7 | oracle* | 98.51 | - | 0.23 | 0.20 | - |
| 8 | perfect* | 95.00 | 95.00 | 0.22 | 0.20 | - |

Table 3: Average coverage and length of confidence intervals over all the entries and an average value of the tuning parameter $\lambda$. The significance level is 0.05.

Model 3: $\Theta_{ij} = 0.5^{i-j}$

$p = 100, n = 200$

|   | Method | Coverage | Length | Average $\lambda$ |
|---|--------|----------|--------|-----------|
| 1 | glasso | 90.43 | 0.19 | 0.138 |
| 2 | glasso-weigh | 75.81 | 0.33 | 0.138 |
| 3 | node-sqrt | 93.36 | 0.28 | 0.152 |
| 4 | node-sqrt-tau | 92.91 | 0.22 | 0.152 |
| 5 | node | 89.88 | 0.20 | 0.152 |
| 6 | MLE | 80.41 | 0.56 | - |
| 7 | perfect* | 95.00 | 0.28 | - |

# 2 Directed acyclic graphs

In this section, we use the de-biasing ideas to construct confidence intervals for edge weights in directed acyclic graphs (abbreviated as DAGs). A directed acyclic graph is a directed graph (we distinguish between edges $(j, k)$ and $(k, j)$) without directed cycles. We consider the Gaussian DAG model, where the DAG represents the probability distribution of a random vector $(X_1, \ldots, X_p)$ with a Gaussian distribution $\mathcal{N}(0, \Sigma_0)$, where $\Sigma_0 \in \mathbb{R}^{p \times p}$ is an unknown covariance matrix. A Gaussian DAG may be represented by the linear structural equations model

$$X_j = \sum_{k \,\in\, \mathrm{pa}(j)} \beta^0_{k,j} X_k + \epsilon_j, \quad j = 1, \ldots, p,$$

where $\epsilon_1, \ldots, \epsilon_p$ are independent and $\epsilon_j \sim \mathcal{N}(0, (\omega^0_j)^2)$. The set $\mathrm{pa}(j)$ is called the set of parents of a node $j$ and it contains all nodes $k \in \{1, \ldots, p\}$ such that there exists a directed edge $k \to j$.

Our aim is to construct confidence intervals for edge weights $\beta^0_{k,j}$. However, without further conditions, the DAG and the $\beta^0_{k,j}$'s may not be identifiable from the structural equations model. To ensure identifiability, we assume that the error variances are equal: $\omega^0_j = \omega_0$ for all $j = 1, \ldots, p$. We remark that one might equivalently assume that the error variances are all known up to a multiplicative constant. In this setting, the DAG is identifiable as shown in Peters and Bühlmann (2014). Our strategy is to use a two-step procedure: in the first step we use the estimator proposed in van de Geer and Bühlmann (2013) to estimate the ordering of the variables and in the second step, we use a de-biased version of nodewise regression to construct the confidence intervals.

Given an $n \times p$ matrix $X = [X_1, \ldots, X_p]$, with rows being $n$ independent observations from the structural equations model, one may rewrite the above model in a matrix form

$$X = XB_0 + E,$$

where $B_0 := (\beta^0_{k,j})$ is a $p \times p$ matrix with $\beta^0_{j,j} = 0$ for all $j$, and $E$ is an $n \times p$ matrix of noise vectors $E := (\epsilon_1, \ldots, \epsilon_p)$ with columns $\epsilon_j$ independent of $X_k$ whenever $\beta^0_{k,j} \neq 0$. The rows of $E$ are independent $\mathcal{N}(0, \omega^2_0 I)$-distributed random vectors. The model then implies that $X$ has covariance matrix

$$\Sigma_0 = \omega^2_0 ((I - B_0)^{-1})^T (I - B_0)^{-1}.$$

We define the precision matrix (assumed to exist) by $\Theta_0 := \Sigma^{-1}_0$. Notice that

$$\Theta_0 = \frac{1}{\omega^2_0} (I - B_0)(I - B_0)^T.$$

We further consider the class of precision matrices corresponding to DAGs. That is, we let

$$\Theta := \Theta(B, \omega) = \frac{1}{\omega^2}(I - B)(I - B)^T,$$

where $(B, \omega)$ is such that there exists a DAG representing the distribution $\mathcal{N}(0, \Sigma)$ with $\Sigma = \omega^2((I - B)^{-1})^T(I - B)^{-1}$. This means that $\omega > 0$ and

17

$B$ can be written as a lower-diagonal matrix, up to permutation of rows. Further we let $s_B$ denote the number of nonzero entries in $B$, which corresponds to the number of edges in the DAG. Moreover, we denote by $\mathcal{B}$ the set of all edge weights $B$ of DAGs with parameters $(B, \omega)$ which have at most $\alpha n / \log p$ incoming edges (parents) at each node, where $\alpha > 0$ is given.

## 2.1 Maximum likelihood estimator with $\ell_0$-penalization

In the first step, we use an $\ell_0$-penalized maximum likelihood estimator to estimate the DAG. Let $\hat{\Sigma} = X^T X / n$ be the Gram matrix based on the design matrix $X$. The minus log-likelihood is proportional to $\ell(\Theta) = \text{trace}(\Theta \hat{\Sigma}) - \log \det(\Theta)$. Consider the penalized maximum likelihood estimator proposed in van de Geer and Bühlmann (2013),

$$
\begin{aligned}
(\hat{B}, \hat{\omega}) \quad := \quad & \text{argmin}_{B, \omega} \{ \ell(B, \omega) + \lambda^2 s_B : \Theta = \Theta(B, \omega), \text{ for some DAG} \\
& \text{with parameters } (B, \omega) \text{ where } B \in \mathcal{B} \}, \quad\quad (17)
\end{aligned}
$$

where $\lambda \geq 0$ is a tuning parameter. The estimator is denoted by $\hat{\Theta} = \Theta(\hat{B}, \hat{\omega})$ and it has $\hat{s} := s_{\hat{B}}$ edges. Calculating the $\ell_0$-penalized maximum likelihood estimator over the class of DAGs is a computationally intensive task, especially because it involves a search through a class of DAGs under a non-convex constraint of acyclicity of the graph and due to the $\ell_0$-penalty. For large scale problems, greedy algorithms may be used, see e.g. Chickering (2002); Hauser and Bühlmann (2012). The reason for using the $\ell_0$-penalty instead of $\ell_1$-penalization in the definition of (17) was discussed in van de Geer and Bühlmann (2013). The $\ell_1$-penalty leads to an objective function which is not constant over equivalent DAGs encoding the same distribution. The $\ell_0$-penalization leads to invariant scores over equivalent DAGs. The theoretical properties of $\hat{\Theta}$ were studied in van de Geer and Bühlmann (2013) under the conditions summarized below. We remark that the paper van de Geer and Bühlmann (2013) primarily studies the estimator (17) with unequal variances and shows that the estimator converges to some member of the Markov equivalence class (cf. Pearl (2016)) of a DAG with a minimal number of edges, under certain conditions.

To make their result precise, we define some further notions. For any vector $\beta \in \mathbb{R}^p$, let $\|X\beta\| := (\beta^T \Sigma_0 \beta)^{1/2}$. By an ordering of variables we mean any permutation of the set $\{1, \ldots, p\}$. For any ordering of the variables, $\pi$, we let $\tilde{B}(\pi)$ be the matrix obtained by doing a Gram-Schmidt orthogonalization of the columns of $X$ in the ordering given by $\pi$, with respect to the norm $\|\cdot\|$. Moreover, let $\tilde{\Omega}_0(\pi) = (I - \tilde{B}_0(\pi))^T \Sigma_0 (I - \tilde{B}_0(\pi)) = \text{diag}((\tilde{\omega}_1^0(\pi))^2, \ldots, (\tilde{\omega}_p^0(\pi))^2)$. We restate the conditions assumed in van de Geer and Bühlmann (2013).

**Condition B1.** *There exists a universal constant $L \geq 1$ such that*

$$
1/L \leq \Lambda_{\min}(\Sigma_0) \leq \Lambda_{\max}(\Sigma_0) \leq L.
$$

**Condition B2.** *There exists a constant $\eta_\omega > 0$ such that for all $\pi$ such that $\tilde{\Omega}_0(\pi) \neq \omega_0^2 I$ it holds*

$$\frac{1}{p}\sum_{i=1}^{p}(|\tilde{\omega}_j^0(\pi)|^2 - \omega_0^2)^2 > 1/\eta_\omega.$$

**Condition B3.** *There exists a sufficiently small constant $\alpha_*$ such that $p \leq \alpha_* n/\log p$.*

Condition B2 is an "omega-min" condition: it imposes that if one uses the wrong permutation then the error variances are far enough from being equal.

Under the above conditions, the $\ell_0$-penalized maximum likelihood estimator with high-probability correctly estimates the ordering of the variables as shown in van de Geer and Bühlmann (2013). Let $\pi_0$ be an ordering of the variables such that a Gram-Schmidt orthogonalization of the columns of $X$ in the order given by $\pi_0$ with respect to the norm $\|\cdot\|$ yields $B_0$. Denote the ordering of variables estimated by the $\ell_0$-penalized maximum likelihood estimator by $\hat{\pi}$. Then the result in van de Geer and Bühlmann (2013) states that under Conditions B1, B2 and B3, with high-probability it holds that $\hat{\pi} = \pi_0$.

## 2.2 Inference for edge weights

Given that we have recovered the true ordering $\pi_0$, the problem reduces to estimation of regression coefficients in a nodewise regression model, where each variable is a function of a known set of its "predecessors". Therefore to construct asymptotically normal estimators for the $\beta_{k,j}^0$'s, we may use a nodewise regression approach.

An estimated ordering $\hat{\pi}$ yields estimates $\hat{\mathrm{p}}(j)$ of the predecessor sets for each node $j = 1, \ldots, p$. If we have recovered the true ordering, that is $\hat{\pi} = \pi_0$, the estimated predecessor sets $\hat{\mathrm{p}}(j)$ are equal to the true predecessor sets, which are supersets of the parent sets $\mathrm{pa}(j)$ for each $j = 1, \ldots, p$. Consequently, given the predecessor sets, we may obtain a new estimator for the edge weights by regressing the $j$-th variable $X_j$ on all its predecessors. The predecessor sets $\hat{\mathrm{p}}(j)$ might be as large as $p - 1$, therefore it is necessary to use regularization. Then we use the de-biasing technique in a similar spirit as in Section 1.2. We remark that in the initial step, one may use any estimator which guarantees exact recovery of the ordering $\pi_0$.

For any non-empty subset $T \subseteq \{1, \ldots, p\}$, we denote by $X_T$ the $n \times |T|$ matrix formed by taking the columns $X_k$ of $X$ such that $k \in T$. We define the nodewise regression estimator as proposed in Janková and van de Geer (2016b) (altenatively, one may use the nodewise square-root Lasso studied in Section 1.4)

$$\hat{\beta}_j = \underset{\beta \in \mathbb{R}^{|\hat{\mathrm{p}}(j)|}}{\operatorname{argmin}} \|X_j - X_{\hat{\mathrm{p}}(j)}\beta\|_2^2/n + 2\lambda_j\|\beta\|_1. \tag{18}$$

The Karush-Kuhn-Tucker conditions for the above optimization problem give

$$-X_{\widehat{p}(j)}^T(X_j - X_{\widehat{p}(j)}\hat{\beta}_j)/n + \lambda_j \hat{Z}_j = 0,$$

where the entries of $\hat{Z}_j$ satisfy $\hat{Z}_{k,j} = \mathrm{sign}(\hat{\beta}_{k,j})$ if $\hat{\beta}_{k,j} \neq 0$, and $\|\hat{Z}_j\|_\infty \leq 1$ ($\hat{\beta}_{k,j}$ denotes the $k$-th entry of $\hat{\beta}_j$). Similarly as in the case of undirected graphical models, we can define a de-biased estimator. The Hessian matrix of the risk function in (18) is given by

$$\hat{\Sigma}_{\widehat{p}(j)} := X_{\widehat{p}(j)}^T X_{\widehat{p}(j)}/n.$$

To find a surrogate inverse for $\hat{\Sigma}_{\widehat{p}(j)}$, we construct $\hat{\Theta}_{\widehat{p}(j)}$ using the nodewise Lasso with tuning parameters $\lambda_{k,j}$ for $k \in \widehat{p}(j)$. Using $\hat{\Theta}_{\widehat{p}(j)}$, we define the de-biased estimator

$$\hat{b}_j := \hat{\beta}_j + \hat{\Theta}_{\widehat{p}(j)}^T X_{\widehat{p}(j)}^T(X_j - X_{\widehat{p}(j)}\hat{\beta}_j)/n. \tag{19}$$

Theorem 7 below shows that the entries of the de-biased estimator are asymptotically normal. To formulate the result, we define $\Theta_{p(j)}^0$ to be the matrix obtained by taking the rows and columns of $\Theta_0$ contained in the true predecessor set $p(j)$. Denote the $k$-th column of $\Theta_{p(j)}^0$ by $\Theta_{p(j),k}^0$. To provide asymptotically normal estimators for the parameters $\beta_{p(j)}^0 = (\beta_{k,j}^0 : k \in p(j))$, we need to impose a sparsity condition on the sizes of the parent sets, which will be denoted by $d_j = |\mathrm{pa}(j)|$.

**Theorem 7** (Regime $p \leq n$). *Let $\hat{B}$ be the estimator defined by (17) with $\lambda \asymp \sqrt{\log p/n}$ and denote the predecessor sets estimated based on $\hat{B}$ by $\hat{p}(j)$ for $j = 1, \ldots, p$. Let $\hat{b}_j$ be defined in (19) with sufficiently large tuning parameters $\lambda_j \asymp \lambda_{k,j} \asymp \sqrt{\log p/n}$, uniformly in $j, k$, where $k \in \widehat{p}(j)$. Assume Conditions B1, B2 and B3 are satisfied with $1/(|\alpha_*| + |\eta_\omega|) = \mathcal{O}(1)$ and assume that $d_j = o(\sqrt{n}/\log p)$. Then it holds*

$$\hat{b}_j - \beta_{p(j)}^0 = (\Theta_{p(j)}^0)^T X_{p(j)}^T \epsilon_j/n + \mathrm{rem}, \tag{20}$$

*where*

$$\|\mathrm{rem}\|_\infty = o_P(1/\sqrt{n}).$$

*Furthermore, for every $k \in p(j)$,*

$$\sqrt{n}(\hat{b}_{k,j} - \beta_{k,j}^0)/\sigma_{k,j} \rightsquigarrow \mathcal{N}(0,1),$$

*where the asymptotic variance of the de-sparsified estimator is given by*

$$\sigma_{k,j}^2 := n\mathrm{var}((\Theta_{p(j),k}^0)^T X_{p(j)}^T \epsilon_j) = \omega_0^2(\Theta_{p(j)}^0)_{kk}.$$

The result of Theorem 7 can be used to construct confidence intervals for the edge weights $\beta_{k,j}^0$. To estimate the asymptotic variance, we may define $\hat{\omega}_j^2 := \|X_j - X_{\widehat{p}(j)}\hat{\beta}_j\|_2^2/n$ and $\hat{\sigma}_{k,j}^2 := \hat{\omega}_j^2(\hat{\Theta}_{\widehat{p}(j)})_{kk}$. The consistency of this estimator may be easily checked.

# 3 Conclusion

We have provided a unified approach to construct asymptotically linear and normal estimators of low-dimensional parameters of the precision matrix based on regularized estimators. These estimators allow us to construct confidence intervals for edge weights in high-dimensional Gaussian graphical models and, under an identifiability condition, for edge weights in the high-dimensional Gaussian DAG model.

For Gaussian graphical models, we provided two explicit simple constructions: one based on a global method using the graphical Lasso and the second based on a local method using nodewise Lasso regressions. Efficient computational methods are available for both methods as discussed in Section 1.5. The constructed estimators are asymptotically normal per entry, achieving the efficient asymptotic variance from the parametric setting. For a detailed analysis of semi-parametric efficiency bounds in Gaussian graphical models, we refer to Janková and van de Geer (2016a). For testing hypothesis about a set of edges, the usual multiple testing corrections may be used although in practical applications, these might turn out to be too conservative. More efficient methods for multiple testing in this setting are yet to be developed. While throughout the presented results we have imposed "exact" sparsity constraints on the underlying parameters, we remark that the results might as well be extended to models which are only approximately sparse (see e.g. Bühlmann and van de Geer (2011)).

Our main interest lied in developing methodology for graphical models representing continuous random vectors. However, many applications involve *discrete* graphical models, where random variables $X_j$ at each vertex $j \in \mathcal{V}$ take values in a discrete space. A popular family of distributions for the binary case where $X_j \in \{-1, 1\}$ is the Ising model. This model finds applications in statistical physics, neuroscience or modeling of social networks. The Ising model can be efficiently estimated via a nodewise method: the individual neighbourhoods can be estimated with $\ell_1$-penalized logistic regression as proposed in Ravikumar, Wainwright, Lafferty, et al. (2010). Logistic regression falls into the framework of generalized linear models for which the de-biasing methodology was proposed in van de Geer et al. (2014). Consequently, one may compute the neighbourhood estimator via $\ell_1$-penalized logistic regression and then compute the de-biased estimator along the lines of van de Geer et al. (2014).

For directed acyclic graphs, we showed that confidence intervals for edge weights may be constructed for the Gaussian DAG when it is identifiable and $p \le \alpha_* n / \log p$. To this end, we require that the error variances in the structural equations model are equal, or known up to a multiplicative constant. If the variance are not equal, the model may not be identifiable and work on inference in this setting is yet to be developed.

# 4 Proofs

## 4.1 Proofs for undirected graphical models

**Lemma 2.** *Assume that $1/L \leq \Lambda_{\min}(\Theta_0) \leq \Lambda_{\max}(\Theta_0) \leq L$ for some constant $L \geq 1$. Let $\mathcal{E}(\Delta) := \mathrm{tr}[\Delta\Sigma_0] - [\log\det(\Delta + \Theta_0) - \log\det(\Theta_0)]$. Then for all $\Delta$ such that $\|\Delta\|_F \leq 1/(2L)$, $\mathcal{E}(\Delta)$ is well defined and*

$$\mathcal{E}(\Delta) \geq \frac{1}{2(L + 1/(2L))^2}\|\Delta\|_F^2. \tag{21}$$

*Proof of Lemma 2.* First we show that $\mathcal{E}(\Delta)$ is well defined for all $\Delta$ such that $\|\Delta\|_F \leq 1/(2L)$. To this end, we need to check that $\Lambda_{\min}(\Theta_0 + \Delta) \geq c_1$ for some $c_1 > 0$. Denote the spectral norm of a matrix $M$ by $\|M\| := \sqrt{\Lambda_{\max}(MM^T)}$. We have

$$\Lambda_{\min}(\Theta_0 + \Delta) = \min_{\|x\|_2 = 1} x^T(\Theta_0 + \Delta)x \geq \Lambda_{\min}(\Theta_0) - \|\Delta\|_F \geq 1/(2L),$$

where we used that $|x^T\Delta x| \leq \|\Delta\|x^Tx$ and that $\|\Delta\| \leq \|\Delta\|_F$.

A second order Taylor expansion with remainder in integral form yields

$$\log\det(\Delta + \Theta_0) - \log\det(\Theta_0)$$

$$= \mathrm{tr}(\Delta\Sigma_0) - \mathrm{vec}(\Delta)^T \left(\int_0^1 (1-v)(\Theta_0 + v\Delta)^{-1} \otimes (\Theta_0 + v\Delta)^{-1} dv\right) \mathrm{vec}(\Delta).$$

Then for all $\Delta$ such that $\|\Delta\|_F \leq 1/(2L)$, it holds

$$\mathcal{E}(\Delta) = \mathrm{vec}(\Delta)^T \left(\int_0^1 (1-v)(\Theta_0 + v\Delta)^{-1} \otimes (\Theta_0 + v\Delta)^{-1} dv\right) \mathrm{vec}(\Delta),$$

where $\otimes$ denotes the Kronecker product and the remainder in the Taylor expansion is in the integral form. Using the fact that the eigenvalues of Kronecker product of symmetric matrices is the product of eigenvalues of the factors, it follows for all $\Delta$ such that $\|\Delta\|_F \leq 1/(2L)$

$$\Lambda_{\min}\left(\int_0^1 (1-v)(\Theta_0 + v\Delta)^{-1} \otimes (\Theta_0 + v\Delta)^{-1} dv\right)$$

$$\geq \int_0^1 (1-v)\Lambda_{\min}^2((\Theta_0 + v\Delta)^{-1}) dv$$

$$\geq \frac{1}{2} \min_{0 \leq v \leq 1} \Lambda_{\min}^2((\Theta_0 + v\Delta)^{-1})$$

$$\geq \frac{1}{2} \min_{\Delta : \|\Delta\|_F \leq 1/(2L)} \Lambda_{\min}^2((\Theta_0 + \Delta)^{-1}).$$

Next we obtain

$$\Lambda_{\min}^2((\Theta_0 + \Delta)^{-1}) = \Lambda_{\max}^{-2}(\Theta_0 + \Delta) \geq (\|\Theta_0\| + \|\Delta\|)^{-2} \geq \frac{1}{(L + 1/(2L))^2} > 0,$$

where we used $\|\Delta\| \leq \|\Delta\|_F \leq 1/(2L)$. Finally this yields that $\mathcal{E}(\Delta) \geq \frac{1}{2(L+1/(2L))^2}\|\Delta\|_F^2$ for all $\|\Delta\|_F \leq 1/(2L)$, as required. $\quad\square$

*Proof of Theorems 1 and 2.* We will prove both Theorem 1 and Theorem 2 at the same time, since the proofs only differ slightly. For the proof of Theorem 2, one has to replace $\hat{\Sigma}$, $\Sigma_0$, $\hat{\Theta}$, $\Theta_0$ in the proof below by $\hat{\Gamma}$, $\Gamma_0$, $\hat{\Theta}_{\text{norm}}$, $K_0$, respectively.

Let $\tilde{\Theta} := \alpha\hat{\Theta} + (1-\alpha)\Theta_0$, where $\alpha := \frac{M}{M+\|\hat{\Theta}-\Theta_0\|_F}$, for some $M > 0$ to be specified later. The definition of $\tilde{\Theta}$ implies that $\|\tilde{\Theta} - \Theta_0\|_F \leq M$. By the convexity of the loss function and by the definition of $\hat{\Theta}$, we have

$$\text{tr}(\tilde{\Theta}\hat{\Sigma}) - \log\det(\tilde{\Theta}) + \lambda\|\tilde{\Theta}^-\|_1 \leq \text{tr}(\Theta_0\hat{\Sigma}) - \log\det(\Theta_0) + \lambda\|\Theta_0^-\|_1. \quad (22)$$

Denote $\Delta = \tilde{\Theta} - \Theta_0$ and let

$$\mathcal{E}(\Delta) := \text{tr}(\Delta\Sigma_0) - [\log\det(\Delta + \Theta_0) - \log\det(\Theta_0)].$$

The inequality (22) implies the basic inequality

$$\mathcal{E}(\Delta) + \lambda\|\tilde{\Theta}^-\|_1 \leq -\text{tr}[\Delta(\hat{\Sigma} - \Sigma_0)] + \lambda\|\Theta_0^-\|_1.$$

On the set $\{\|\hat{\Sigma} - \Sigma_0\|_\infty \leq \lambda_0\}$, we can bound the empirical process term by

$$\begin{aligned}
|\text{tr}[\Delta(\hat{\Sigma} - \Sigma_0)]| &\leq& \|\hat{\Sigma}^- - \Sigma_0^-\|_\infty\|\Delta^-\|_1 + \|\hat{\Sigma}^+ - \Sigma_0^+\|_2\|\Delta^+\|_F \\
&\leq& \lambda_0\|\Delta^-\|_1 + \|\hat{\Sigma}^+ - \Sigma_0^+\|_2\|\Delta^+\|_F. \quad (23)
\end{aligned}$$

In what follows, we work on the set $\{\|\hat{\Sigma} - \Sigma_0\|_\infty \leq \lambda_0\}$.
We now choose $M$ such that $M \leq 1/(2L)$, this then implies $\|\tilde{\Theta}-\Theta_0\|_F \leq 1/(2L)$. But then Lemma 2 implies that $\mathcal{E}(\tilde{\Theta} - \Theta_0)$ is well defined and

$$\mathcal{E}(\tilde{\Theta} - \Theta_0) \geq c\|\tilde{\Theta} - \Theta_0\|_F^2, \quad (24)$$

where one can take $c := 1/(8L^2)$. Using bounds (23) and (24), we obtain from the basic inequality

$$c\|\Delta\|_F^2 + \lambda\|\tilde{\Theta}^-\|_1 \leq \lambda_0\|\Delta^-\|_1 + \lambda\|\Theta_0^-\|_1 + \|\hat{\Sigma}^+ - \Sigma_0^+\|_2\|\Delta^+\|_F$$

By the triangle inequality and taking $\lambda \geq 2\lambda_0$, we obtain

$$2c\|\Delta\|_F^2 + \lambda\|\tilde{\Theta}_{S^c}^-\|_1 \leq 3\lambda\|\Delta_S^-\|_1 + 2\|\hat{\Sigma}^+ - \Sigma_0^+\|_2\|\Delta^+\|_F$$

Consequently,

$$\begin{aligned}
2c\|\Delta\|_F^2 + \lambda\|\Delta^-\|_1 &\leq& 4\lambda\|\Delta_S^-\|_1 + 2\|\hat{\Sigma}^+ - \Sigma_0^+\|_2\|\Delta^+\|_F \\
&\leq& 4\lambda\sqrt{s}\|\Delta_S^-\|_F + 2\|\hat{\Sigma}^+ - \Sigma_0^+\|_2\|\Delta^+\|_F \\
&\leq& 8s\lambda^2/c^2 + c\|\Delta_S^-\|_F^2/2 + 8\|\hat{\Sigma}^+ - \Sigma_0^+\|_2^2/c + c\|\Delta^+\|_F^2/2.
\end{aligned}$$

Taking $M$ such that $\lambda_0 M \geq 8s\lambda^2/c^2 + 8\|\hat{\Sigma}^+ - \Sigma_0^+\|_2^2/c$,

$$c\|\Delta\|_F^2 + \lambda\|\Delta^-\|_1 \leq 8s\lambda^2/c^2 + 8\|\hat{\Sigma}^+ - \Sigma_0^+\|_2^2/c \leq \lambda_0 M.$$

Taking $M \geq 4\lambda_0/c$,

$$\|\Delta\|_F^2 \leq \lambda_0 M/c \leq M^2/4.$$

But then $\|\Delta\|_F \leq M/2$. The definition of $\tilde{\Theta}$ in turn implies that $\|\hat{\Theta}-\Theta_0\|_F \leq M$, and we can repeat all the arguments with $\hat{\Theta}$ in place of $\tilde{\Theta}$. Repetition of the arguments leads to the oracle inequality

$$c\|\hat{\Theta}-\Theta_0\|_F^2 + \lambda\|\hat{\Theta}^- - \Theta_0^-\|_1 \leq 8s\lambda^2/c^2 + 8\|\hat{\Sigma}^+ - \Sigma_0^+\|_2^2/c.$$

Finally we distinguish the case of non-normalized graphical Lasso (based on the covariance matrix) and the normalized graphical Lasso (based on the correlation matrix). We have for the case of

a) <u>normalized graphical Lasso</u>: $\hat{\Sigma}^+ - \Sigma_0^+ = 0$ (recall here that $\hat{\Sigma} \equiv \hat{R}, \Sigma_0 \equiv R_0$) and the oracle inequality gives

$$c\|\hat{\Theta}-\Theta_0\|_F^2 + \lambda\|\hat{\Theta}^- - \Theta_0^-\|_1 \leq 8s\lambda^2/c^2.$$

b) <u>non-normalized graphical Lasso</u>: we can bound

$$\|\hat{\Sigma}^+ - \Sigma_0^+\|_2 \leq \sqrt{p}\|\hat{\Sigma}^+ - \Sigma_0^+\|_\infty \leq \sqrt{p}\lambda_0.$$

Hence the oracle inequality gives

$$c\|\hat{\Theta}-\Theta_0\|_F^2 + \lambda\|\hat{\Theta}^- - \Theta_0^-\|_1 \leq 8s\lambda^2/c^2 + 8p\lambda_0^2/c.$$

To show the second statement of the theorems, we use the above oracle inequalities and the following upper bound

$$\begin{aligned}
\left\|\left\|\hat{\Theta}-\Theta_0\right\|\right\|_1 &\leq \max_{j=1,\ldots,p}|\hat{\Theta}_{jj} - \Theta_{jj}^0| + \|\hat{\Theta}_j^- - (\Theta_j^0)^-\|_1 \\
&\leq \|\hat{\Theta}-\Theta_0\|_F + \|\hat{\Theta}^- - \Theta_0^-\|_1.
\end{aligned}$$

To show the third statement of Theorem 2, we use the upper bound

$$\begin{aligned}
\left\|\left\|\hat{\Theta}_{\mathrm{w}}-\Theta_0\right\|\right\|_1 &= \left\|\left\|\hat{W}^{-1}\hat{\Theta}_{\mathrm{norm}}\hat{W}^{-1} - W_0^{-1}K_0W_0^{-1}\right\|\right\|_1 \\
&\leq \|\hat{W}\|_\infty^2\left\|\left\|\hat{\Theta}_{\mathrm{norm}} - K_0\right\|\right\|_1 + \|\hat{W}-W_0\|_\infty\|\|K_0\|\|_1\|\hat{W}\|_\infty \\
&\quad + \|W_0\|_\infty\|\|K_0\|\|_1\|\hat{W}-W_0\|_\infty
\end{aligned}$$

$\square$

*Proof of Theorem 3.* Denote $C_L := 16(8L^2)^2$. Using the results of Theorem 1, we obtain

$$\begin{aligned}
\|\mathrm{rem}\|_\infty &\leq \|(\hat{\Theta}-\Theta_0)^T(\hat{\Sigma}\Theta_0 - I)\|_\infty + \|(\hat{\Theta}-\Theta_0)^T\lambda\hat{Z}\hat{\Theta}\|_\infty \\
&\leq \left\|\left\|\hat{\Theta}-\Theta_0\right\|\right\|_1\|\hat{\Sigma}-\Sigma_0\|_\infty\|\|\Theta_0\|\|_1 + \left\|\left\|\hat{\Theta}-\Theta_0\right\|\right\|_1\lambda\|\hat{Z}\|_\infty\left\|\left\|\hat{\Theta}\right\|\right\|_1 \\
&\leq C_L(p+s)\lambda\sqrt{d+1}\Lambda_{\max}(\Theta_0)\lambda_0 + 2C_L(p+s)\lambda\sqrt{d+1}\Lambda_{\max}(\Theta_0)\lambda \\
&\leq \frac{3}{2}C_L L(p+s)\sqrt{d+1}\lambda^2.
\end{aligned}$$

24

Taking $\lambda \asymp \sqrt{\log p/n}$, Lemma 1 implies $\|\hat{\Sigma} - \Sigma_0\|_\infty = \mathcal{O}_P(\sqrt{\log p/n})$. Then by the sparsity condition, we obtain $\|\mathrm{rem}\|_\infty = \mathcal{O}_P(1/\sqrt{n})$. By Conditions A1 and A2, the random variable $(\Theta_0(\hat{\Sigma} - \Sigma_0)\Theta_0)_{ij}$ has bounded fourth moments and asymptotic normality per entry follows by application of Lindeberg's central limit theorem for triangular arrays (see Janková and van de Geer (2016b) for more details). $\qquad\square$

*Proof of Theorem 4.* For the remainder, we obtain similarly as in the proof of Theorem 1, $\|\mathrm{rem}\|_\infty \le \frac{3}{2}C_L L s\sqrt{d+1}\lambda^2$. Asymptotic normality follows by analogous arguments. $\qquad\square$

*Proof of Proposition 1.* Denote $C_L := 16(8L^2)^2$. Using the results of Theorem 2, we obtain

$$
\begin{aligned}
\|\mathrm{rem}\|_\infty &\le \left\|\left\|\left\|\hat{\Theta}_{\mathrm{norm}} - K_0\right\|\right\|\right\|_1 \|\hat{\Gamma}K_0 - I\|_\infty + \left\|\left\|\left\|\hat{\Theta}_{\mathrm{norm}} - K_0\right\|\right\|\right\|_1 \lambda \|\hat{Z}\|_\infty \left\|\left\|\left\|\hat{\Theta}_{\mathrm{norm}}\right\|\right\|\right\|_1 \\
&\le C_L s\lambda\sqrt{d+1}\Lambda_{\max}(\Theta_0)\lambda_0 + C_L s\lambda\sqrt{d+1}\Lambda_{\max}(\Theta_0)\lambda \\
&\le \frac{3}{2}C_L L s\sqrt{d+1}\lambda^2.
\end{aligned}
$$

The sparsity condition implies the result. $\qquad\square$

*Proof of Theorem 6.* The proof follows along the same lines as the proof of Theorem 1 in Janková and van de Geer (2016b). The only difference is that here we consider a weighted Lasso to estimate the partial correlations and the estimator $\hat{\tau}_j$ is defined slightly differently. But for the weighted Lasso (with weights bounded away from zero and bounded from above with high probability), oracle inequalities of the same order can be obtained, see Section 6.9 in Bühlmann and van de Geer (2011), i.e.

$$
\|X_{-j}(\hat{\gamma}_j - \gamma_j^0)\|_2^2/n + \lambda\|\hat{\gamma}_j - \gamma_j^0\|_1 = \mathcal{O}_P(d_j \log p/n).
$$

For the estimator of variance we have

$$
\begin{aligned}
|\hat{\tau}_j^2 - \tau_j^2| &\le \|X_{-j}(\hat{\gamma}_j - \gamma_j^0)\|_2^2/n + 2|(X_j - X_{-j}\gamma_j^0)^T X_{-j}(\hat{\gamma}_j - \gamma_j^0)/n| \\
&= \|X_{-j}(\hat{\gamma}_j - \gamma_j^0)\|_2^2/n + 2\|(X_j - X_{-j}\gamma_j^0)^T X_{-j}/n\|_\infty\|\hat{\gamma}_j - \gamma_j^0\|_1 \\
&= \mathcal{O}_P(1/\sqrt{n}).
\end{aligned}
$$

The rest of the proof follows as in Janková and van de Geer (2016b).

$\qquad\square$

## 4.2 Proofs for directed acyclic graphs

*Proof of Theorem 7.* By Theorem 5.1 in van de Geer and Bühlmann (2013), we have under the conditions of the theorem that $\hat{\pi} = \pi_0$ with high probability. Then also $\hat{\mathrm{p}}(j) = \mathrm{p}(j)$ for all $j$, with high probability. Therefore, the estimated $\hat{\mathrm{p}}(j)$ in the definitions of $\hat{\beta}_j$ and $\hat{\Theta}_{\hat{\mathrm{p}}(j),k}, k \in \hat{\mathrm{p}}(j)$ (and elsewhere) can be replaced

25

by $\mathrm{p}(j)$. The nodewise Lasso then yields oracle estimators $\hat{\beta}_j$ and $\hat{\Theta}_k, k \in \mathrm{p}(j)$ under the Condition B1 and under the sparsity $d_j = o(\sqrt{n/\log p})$ (see Janková and van de Geer (2016b)). This gives in particular that for all $j = 1, \ldots, p$

$$\max_k \|\hat{\Theta}_{\mathrm{p}(j),k} - \Theta^0_{\mathrm{p}(j),k}\|_1 = \mathcal{O}_P(\max_k d_j \lambda_j),$$

$$\max_k \|\hat{\Sigma}_{\mathrm{p}(j)} \hat{\Theta}_{\mathrm{p}(j),k} - e_k\|_\infty = \mathcal{O}_P(\max_k \lambda_j),$$

$$\|\hat{\beta}_j - \beta^0_j\|_1 = \mathcal{O}_P(\max_{j=1,\ldots,p} d_j \lambda_j).$$

We can write the decomposition

$$\hat{b}_{k,j} - \beta^0_{k,j} = (\Theta^0_{\mathrm{p}(j),k})^T X^T_{\mathrm{p}(j)} \epsilon_j / n + \mathrm{rem}_{k,j}, \tag{25}$$

where $\mathrm{rem}_{k,j} = (\hat{\Theta}_{\mathrm{p}(j),k} - \Theta^0_{\mathrm{p}(j),k})^T X^T_{\mathrm{p}(j)} \epsilon_j / n - (\hat{\Sigma}_{\mathrm{p}(j)} \hat{\Theta}_{\mathrm{p}(j),k} - e_k)^T (\hat{\beta}_j - \beta^0_j)$. First note that by normality and by the independence of $X_{\mathrm{p}(j)}$ and $\epsilon_j$ (which follows by the independence of $\epsilon_j$'s and acyclicity of the graph), it holds $\|X^T_{\mathrm{p}(j)} \epsilon_j / n\|_\infty = \mathcal{O}_P(\sqrt{\log p/n})$. By Hölder's inequality

$$\begin{aligned}
\max_k |\mathrm{rem}_{k,j}| &\leq \max_k \|\hat{\Theta}_{\mathrm{p}(j),k} - \Theta^0_{\mathrm{p}(j),k}\|_1 \|X^T_{\mathrm{p}(j)} \epsilon_j / n\|_\infty \\
&+ \max_k \|\hat{\Theta}^T_{\mathrm{p}(j),k} \hat{\Sigma}_{\mathrm{p}(j)} - e_k\|_\infty \|\hat{\beta}_j - \beta^0_j\|_1 \\
&= \mathcal{O}_P(d_j \log p/n) = o_P(1/\sqrt{n}),
\end{aligned}$$

where we used the sparsity assumption $d_j = o(\sqrt{n}/\log p)$. Thus we have shown that the remainder in (25) is of small order $1/\sqrt{n}$. Then applying Lindeberg's central limit theorem for triangular arrays and by Conditions A1 and A2,

$$(\Theta^0_{\mathrm{p}(j),k})^T X^T_{\mathrm{p}(j)} \epsilon_j / (\sigma_{k,j} \sqrt{n}) \rightsquigarrow \mathcal{N}(0,1),$$

which shows the claim. $\qquad\square$

# References

A. Belloni, V. Chernozhukov, and L. Wang. Square-Root Lasso: Pivotal Recovery of Sparse Signals via Conic Programming. *Biometrika*, 98(4):791–806, 2011.

P. J. Bickel and E. Levina. Regularized estimation of large covariance matrices. *Annals of Statistics*, 36(1):199–227, 2008a.

P. J. Bickel and E. Levina. Covariance regularization by thresholding. *Annals of Statistics*, 36(6):2577–2604, 2008b.

P. Bühlmann and S. van de Geer. Statistics for high-dimensional data. *Springer*, 2011.

T. Cai, W. Liu, and X. Luo. A constrained l1 minimization approach to sparse precision matrix estimation. *Journal of the American Statistical Association*, 106, 2011.

E. Candes and T. Tao. The dantzig selector: Statistical estimation when p is much larger than n. *Annals of Statistics*, 35(6):2313–2351, 12 2007. doi: 10.1214/009053606000001523. URL http://dx.doi.org/10.1214/009053606000001523.

V. Chernozhukov, D. Chetverikov, and K. Kato. Central Limit Theorems and Bootstrap in High Dimensions. *ArXiv: 1412.3661*, 2014.

D.M. Chickering. Optimal structure identification with greedy search. *Journal of Machine Learning Research*, 3:507–554, 2002.

Alexandre d'Aspremont, Onureena Banerjee, and Laurent El Ghaoui. First-order methods for sparse covariance selection. *SIAM Journal on Matrix Analysis and Applications*, 30(1):56–66, 2008. ISSN 0895-4798. doi: 10.1137/060670985. URL http://dx.doi.org/10.1137/060670985.

Mathias Drton and Marloes H Maathuis. Structure learning in graphical modeling. *Annual Review of Statistics and Its Application*, 4:365–393.

B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least Angle Regression. *Ann. Statist.*, 32(2):407–451, June 2004.

N. El Karoui. Operator norm consistent estimation of large dimensional sparse covariance matrices. *Annals of Statistics*, 36(6):2717–2756, 2008.

J. Friedman, T. Hastie, and R. Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9:432–441, 2008.

A. Hauser and P. Bühlmann. Characterization and greedy learning of interventional Markov equivalence classes of directed acyclic graphs. *Journal of Machine Learning Research*, 13:2409–2464, 2012.

J. Janková and S. van de Geer. Confidence intervals for high-dimensional inverse covariance estimation. *Electronic Journal of Statistics*, 9(1):1205 –1229, 2014.

J. Janková and S. van de Geer. Semi-parametric efficiency bounds for high-dimensional models. *ArXiv:1601.00815*, 2016a.

J. Janková and S. van de Geer. Honest confidence regions and optimality for high-dimensional precision matrix estimation. *TEST*, 26(1):143–162, 2016b.

A. Javanmard and A. Montanari. Confidence intervals and hypothesis testing for high-dimensional regression. *Journal of Machine Learning Research*, 15 (1):2869–2909, 2014.

I. M. Johnstone and A. Y. Lu. On consistency and sparsity for principal components analysis in high dimensions. *Journal of the American Statistical Association*, 104(486):682–693, 2009.

Iain M. Johnstone. On the distribution of the largest eigenvalue in principal components analysis. *Annals of Statistics*, 29(2):295–327, 04 2001. doi: 10.1214/aos/1009210544.

Weidong Liu et al. Gaussian graphical model estimation with false discovery rate control. *Annals of Statistics*, 41(6):2948–2978, 2013.

R. Mazumder and T. Hastie. The Graphical Lasso: New Insights and Alternatives. *Electronic Journal of Statistics*, 6:2125–2149, 2012.

N. Meinshausen and P. Bühlmann. High-dimensional graphs and variable selection with the lasso. *Annals of Statistics*, 34(3):1436–1462, 06 2006. doi: 10.1214/009053606000000281. URL http://dx.doi.org/10.1214/009053606000000281.

J. Pearl. *Causality: Models, Reasoning and Inference*. Cambridge University Press, 2016.

J. Peters and P. Bühlmann. Identifiability of Gaussian structural equation models with equal error variances. *Biometrika*, 101:219–228, 2014.

P. Ravikumar, G. Raskutti, M. J. Wainwright, and B. Yu. High-dimensional covariance estimation by minimizing l1-penalized log-determinant divergence. *Electronic Journal of Statistics*, 5:935–980, 2008.

Pradeep Ravikumar, Martin J Wainwright, John D Lafferty, et al. High-dimensional Ising model selection using 1-regularized logistic regression. *Annals of Statistics*, 38(3):1287–1319, 2010.

Zhao Ren, Tingni Sun, Cun-Hui Zhang, Harrison H Zhou, et al. Asymptotic normality and optimalities in estimation of large gaussian graphical models. *Annals of Statistics*, 43(3):991–1026, 2015.

A. J. Rothman, P. J. Bickel, E. Levina, and J. Zhu. Sparse permutation invariant covariance estimation. *Electronic Journal of Statistics*, 2:494–515, 2008. doi: 10.1214/08-EJS176. URL http://dx.doi.org/10.1214/08-EJS176.

T. Sun and C.-H. Zhang. Sparse matrix inversion with scaled Lasso. *Journal of Machine Learning Research*, 14:3385–3418, 2012.

R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of The Royal Statistical Society: Series B*, 58:267–288, 1996.

S. van de Geer. *Estimation and Testing under Sparsity: École d'Été de Saint-Flour XLV*. Springer, 2016.

S. van de Geer and P. Bühlmann. On the conditions used to prove oracle results for the lasso. *Electronic Journal of Statistics*, 3:1360–1392, 2009. doi: 10.1214/09-EJS506. URL http://dx.doi.org/10.1214/09-EJS506.

S. van de Geer and P. Bühlmann. $\ell_0$-penalized maximum likelihood for sparse directed acyclic graphs. *Annals of Statistics*, 41(2):536–567, 2013.

S. van de Geer, P. Bühlmann, Y. Ritov, and R. Dezeure. On asymptotically optimal confidence regions and tests for high-dimensional models. *Annals of Statistics*, 42(3):1166–1202, 2014.

Sara van de Geer. Worst possible sub-directions in high-dimensional models. *Journal of Multivariate Analysis*, 146:248–260, 2016.

Larry Wasserman, Mladen Kolar, Alessandro Rinaldo, et al. Berry-Esseen bounds for estimating undirected graphs. *Electronic Journal of Statistics*, 8(1):1188–1224, 2014.

Ming Yu, Mladen Kolar, and Varun Gupta. Statistical inference for pairwise graphical models using score matching. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 2829–2837. Curran Associates, Inc., 2016.

Ming Yuan. High dimensional inverse covariance matrix estimation via linear programming. *The Journal of Machine Learning Research*, 11:2261–2286, 2010.

Ming Yuan and Yi Lin. Model selection and estimation in the Gaussian graphical model. *Biometrika*, 94(1):19–35, 2007.

C.-H. Zhang and S. S. Zhang. Confidence intervals for low-dimensional parameters in high-dimensional linear models. *Journal of the Royal Statistical Society: Series B*, 76:217–242, 2014.