

Discussion Papers In Economics And Business

Estimation of High Dimensional Vector
Autoregression via Sparse Precision Matrix

Benjamin Poignard

Manabu Asai

Discussion Paper 21-03

April 2021

Graduate School of Economics
Osaka University, Toyonaka, Osaka 560-0043, JAPAN

Estimation of High Dimensional Vector Autoregression via Sparse Precision Matrix*

Benjamin Poignard[†]

Graduate School of Economics, Osaka University and Riken-AIP, Japan

Manabu Asai[‡]

Faculty of Economics, Soka University, Japan

Abstract

We consider the problem of estimating sparse structural vector autoregression (SVAR) processes via penalized precision matrix. Such matrix is the output of the underlying directed acyclic graph of the SVAR process, whose zero components correspond to zero SVAR coefficients. The precision matrix estimators are deduced from the class of Bregman divergences and regularized by the SCAD, MCP and LASSO penalties. Under suitable regularity conditions, we derive error bounds for the regularized precision matrix for each Bregman divergence. Moreover, we establish the support recovery property, including the case when the penalty is non-convex. These theoretical results are supported by empirical studies.

Keywords: sparse structural vector autoregression; statistical consistency; support recovery.

*The authors are most grateful to Yoshihisa Baba for very helpful comments and suggestions. The authors acknowledge the financial support of the Japan Society for the Promotion of Science.

[†]E-mail address: bpoignard@econ.osaka-u.ac.jp (corresponding author)

[‡]E-mail address: m-asai@soka.ac.jp

1 Introduction

The Vector AutoRegressive (VAR) process is a standard model for multivariate time series. From an inferential viewpoint, the ordinary least squares method is inapplicable once fitted to high-dimensional data due to a $O(pN^2)$ complexity order, where N is the number of variables and p the number of lags. Numerous studies modelled sparsity among the VAR parameters to tackle the over-fitting issue, which justified the use of penalized OLS losses: see, e.g., Basu and Michailidis (2015), who considered a penalized OLS loss for sparse stable Gaussian VAR processes; in the same manner, Wong, Li and Tewari (2020) derived some consistency results of penalized OLS based estimators for α -mixing Gaussian VAR processes. Rather than specifying sparse VAR matrix parameters, Alquier, Bertin, Doukhan and Garnier (2020) specified a low-rank constraint on the VAR transition matrix in the presence of factors, where the motivation is to improve the prediction accuracy.

Furthermore, the Structural VAR (SVAR) model has been used to accommodate economic theory within the VAR framework: see, e.g., Blanchard and Quah (1989) and Waggoner and Zha (2003) among others. Apart from economic theory, Tunnicliffe-Wilson and Reale (2008), Oxley, Reale, Tunnicliffe-Wilson (2009), and Ahelegbey, Billio and Casarin (2016) discussed the interpretation and identification of SVAR models in terms of graphical modelling. In this paper, we consider the structure of SVAR models and their identification based on a graphical representation, extending Sims (1980) and the previous papers as an approach which is free from the constraints usually assumed in economic theory. We propose a different procedure for the analysis of the issues related to SVAR inference in high dimension: our analysis focuses on the relationships between the SVAR parameters and the precision matrix - i.e., the inverse variance-covariance - of the SVAR process, which allows for a sparse estimation of the SVAR parameters.

More precisely, in this study, we consider the following problem: given T observations of a N -dimensional SVAR(p) process (\mathbf{Y}_t) , estimate the sparse precision matrix Θ of $\mathbf{X}_t = (\mathbf{Y}_{t-p}^\top, \dots, \mathbf{Y}_t^\top)^\top$, whose components provide sparse consistent estimators of the SVAR parameters by the one-to-one matching between the

zero coefficients of the SVAR parameters and the precision matrix. Such matching between the precision matrix Θ and the SVAR parameters relies on the so-called directed acyclic graph, where the elements of \mathbf{X}_t form the nodes and past dependence is captured by the directed edges linking the nodes. By such relationship, we can show that the sparsity on Θ gives a one-to-one correspondence to the zero coefficients of the SVAR(p) parameters. Moreover, our method ensures the stationarity of the resulting process while fostering sparsity on the SVAR coefficients. There exist several methods fostering sparsity for the precision matrix: for example, Bickel and Levina (2008a) considered a thresholding operator for the inverse, which constrains the parameters toward zero should they exceed a certain threshold. Rothman, Levina, and Zhu (2009) showed that this estimator satisfies the “sparsistency” property in the sense of Lam and Fan (2007), that is the true zero parameters are correctly identified with probability tending to one and the sparse estimator is sign consistent for nonzero elements. Motivated by the structure of some time series data, Furrer and Bengtsson (2012) focused on banded precision matrices and fostered sparsity by shrinking the off-diagonal entries based on their distances with respect to the diagonal. Wu and Pourahmadi (2003) and Huang, Liu, Pourahmadi and Liu (2006) considered a Cholesky decomposition of the precision matrix with a thresholding procedure and LASSO penalization of the Cholesky factors respectively. For the aspect of “regression” in the VAR, we may carry out the latter approach based on the relationship between regression coefficients and the precision matrix for vector of relevant variables. The Cholesky decomposition automatically ensures the positive definiteness of the estimated covariance matrix. Bickel and Levina (2008b) provided the conditions for the statistical consistency of such estimators while highlighting that a key feature of such decomposition is that the latter depend on the variable ordering and thus the sparse structure highly depends on such order. Another parsimonious approach consists in assuming a particular sparse structure on the variance covariance or the precision matrix. In particular, Bickel and Levina (2008a) applied a thresholding procedure under the sparse row assumption - the ℓ_q -sparse assumption -. Finally, a significant literature is dedicated to the parsimonious analysis of precision matrices under an element-wise sparsity assumption. Yuan and Lin (2007) or Ravikumar, Wainwright, Raskutti and Yu (2011) assumed an

element-wise sparse precision matrix and proposed a LASSO Gaussian likelihood estimator. Zhang and Zou (2014) proposed an alternative LASSO penalized loss, namely the D-trace loss, to estimate such matrix. In the same vein, Loh and Wainwright (2017) extended the penalization of the Gaussian precision matrix to non-convex penalties, which allows for relaxing the incoherence - or irrepresentability - condition.

Our setting lies within such sparse element-wise precision matrix for estimating sparse SVAR parameters. Our main contributions are as follows: first, we provide a novel estimation method for SVAR parameters under the sparsity assumption; we provide error bounds for a broad range of sparse estimators of Θ in the ℓ_1 , ℓ_2 and ℓ_∞ senses for specific scaling behaviours of (T, d, k_0) , where d and k_0 respectively are the dimension problem and the cardinality of the true unknown sparse support; finally, we provide the conditions to satisfy the support recovery property. The estimators of Θ are deduced from the class of penalized Bregman divergence losses, including the Gaussian, least squares and von Neumann losses, thus providing new estimators for sparse precision matrices. To the best of our knowledge, this paper is the first attempt to link general penalized - potentially non-convex - M-estimators and the Bregman divergence-based inference for sparse precision matrix. Our study shares a similar spirit to that of Ravikumar et al. (2011) and Loh and Wainwright (2017), who derived the conditions for statistical consistency of component-wise sparse Gaussian based precision matrix. But our work differs from these studies in two main respects: we provide bounds on ℓ_1 -, ℓ_2 - and ℓ_∞ - errors and the conditions to satisfy the support recovery property within the framework of Bregman divergence, thus providing new sparse estimators for precision matrix.

The framework we use to derive such error bounds is related to the study of Poignard and Fermanian (2021), which covers a broad range of non-convex objective functions for sparse M-estimation. Assuming the restricted strong convexity (see e.g. Negahban, Ravikumar, Wainwright and Yu, 2012) of the non-penalized loss function and for suitable regularity conditions of the penalty, they derived some error bounds for the penalized estimators. To establish the conditions for support recovery, our proof techniques are inspired from the primal dual witness method of Loh and Wainwright (2017). In our study, we extend these results to a broad range of sparse precision matrix estimators. We quantify the statistical accuracy and discuss the

relevance of these theoretical bounds for each M-criterion. The scaling behaviours with respect to (T, d, k_0) that we derive for support recovery highly depend on the regularity of the Bregman divergence.

The organization of the paper is as follows. In Section 2, we describe our approach for sparse SVAR modelling based on sparse precision matrix. In Section 3, we provide our sparse estimation framework based on Bregman divergences together with the theoretical properties - error bounds and support recovery - of the corresponding estimators. Section 4 illustrates these theoretical properties through simulation and real data experiments. All intermediary results, proofs and figures are contained in Section 5.

Notations. Throughout this paper, we denote the cardinality of a set E by $|E|$. For a vector $\mathbf{v} \in \mathbb{R}^d$, the ℓ_p norm is $\|\mathbf{v}\|_p = (\sum_{k=1}^d |\mathbf{v}_k|^p)^{1/p}$ for $p > 0$, and $\|\mathbf{v}\|_\infty = \max_i |\mathbf{v}_i|$. Let the subset $\mathcal{A} \subseteq \{1, \dots, d\}$, then $\mathbf{v}_{\mathcal{A}} \in \mathbb{R}^{|\mathcal{A}|}$ is the vector \mathbf{v} restricted to \mathcal{A} . We write $\mathcal{M}_{d_1 \times d_2}(\mathbb{R})$ the set of $d_1 \times d_2$ -dimensional matrices with real coefficients. For a matrix A , $\|A\|_s$, $\|A\|_\infty$ and $\|A\|_F$ are the spectral, infinity and Frobenius norms, respectively, and $\|A\|_{\max} = \max_{i,j} |A_{i,j}|$ is the coordinate-wise maximum (in absolute value). We write A^\top (resp. \mathbf{v}^\top) to denote the transpose of the matrix A (resp. the vector \mathbf{v}). We write $\text{vec}(A)$ to denote the vectorization operator that stacks the columns of A on top of one another into a vector. We denote by $A \succ 0$ (resp. $A \succeq 0$) the positive definiteness (resp. semi-definiteness) of A and $\text{vech}(A)$ the $d(d+1)/2$ vector that stacks the columns of the lower triangular part of $A \in \mathcal{M}_{d \times d}(\mathbb{R})$. $\lambda_{\min}(A)$ (resp. $\lambda_{\max}(A)$) denotes the minimum (resp. maximum) eigenvalue of A . For a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, we denote by ∇f the gradient or subgradient of f and $\nabla^2 f$ the Hessian of f . We denote by $(\nabla^2 f)_{\mathcal{A}\mathcal{A}}$ the Hessian of f restricted to the block \mathcal{A} . We write \mathcal{A}^c to denote the complement of the set \mathcal{A} . The expression *with high probability* refers to event occurring with probability approaching one when (T, d, k_0) tend to infinity. The scaling results for (T, d, k_0) are expressed as $f(T) \geq Mg(k_0, d)$ for $0 < M < \infty$ some universal constant and continuous functions $f(\cdot), g(\cdot)$.

2 High Dimensional VAR

Let (\mathbf{Y}_t) be an N -dimensional random vector of time series and consider the following VAR model:

$$\mathbf{Y}_t = A_1 \mathbf{Y}_{t-1} + \cdots + A_p \mathbf{Y}_{t-p} + \mathbf{u}_t, \quad t = 1, \dots, T, \quad (1)$$

where $\mathbf{u}_t \in \mathbb{R}^N$ with $\mathbb{E}[\mathbf{u}_t] = \mathbf{0}$ and $\text{Var}(\mathbf{u}_t) = \Sigma_u \in \mathcal{M}_{N \times N}(\mathbb{R})$, where $\Sigma_u \succ 0$ and for each $i = 1, \dots, p$, $A_i \in \mathcal{M}_{N \times N}(\mathbb{R})$. The driving parameters are A_i, Σ_u . Although we assume that \mathbf{Y}_t is a mean zero vector, our setting can straightforwardly include a constant term in (1). Now consider an invertible matrix $B_0 \in \mathcal{M}_{N \times N}(\mathbb{R})$ with diagonal elements being equal to one. Multiplying B_0 from the left of both sides of (1), we obtain the following form:

$$B_0 \mathbf{Y}_t = B_1 \mathbf{Y}_{t-1} + \cdots + B_p \mathbf{Y}_{t-p} + \mathbf{e}_t, \quad t = 1, \dots, T, \quad (2)$$

where $\mathbf{e}_t = B_0 \mathbf{u}_t$ and $B_i = B_0 A_i$. By construction, $\mathbb{E}[\mathbf{e}_t] = \mathbf{0}$, $\Sigma_e = \text{Var}(\mathbf{e}_t) = B_0 \Sigma_u B_0'$, and Σ_e is positive definite. Equation (2) is known as the ‘A-model’ in the literature on structural VAR (SVAR) models: see, e.g., Lütkepohl (2006, 2017). A challenging task concerns SVAR identification. One way is to assume (i) Σ_e is diagonal, (ii) the diagonal elements of B_0 are one, and (iii) the number of zeros in B_0 is $N(N-1)/2$: these identification restrictions are provided in Proposition 9.1 of Lütkepohl (2006). As pointed out in Subsection 9.1.4 of Lütkepohl (2006), assuming such restrictions is not necessary, especially due to the lack of meaningful economic justifications to impose contemporaneous restrictions. For instance, Blanchard and Quah (1989) developed a framework for imposing long-run restrictions for structural VAR models: such restrictions are based on long-run neutrality properties.

Rather than imposing the above restrictions a priori on B_0 or on long-run representations, we pursue the direction of Sims (1980) for a data-oriented identification based on a graphical representation, which will be explained below. We highlight the existence of an alternative SVAR identification method based on sign restrictions: it consists in dropping doubtful restrictions one after one to identify the most likely admissible model within the set of structural VAR models that satisfy the assumed sign restrictions. For instance,

within the ‘B-model’ described in Subsection 9.1.2 of Lütkepohl (2006), Inoue and Killian (2013) developed an approach based on all admissible models which satisfy signs of parameters derived by economic theory.

Since our analysis focuses on potentially large dimensional VAR model, we assume sparsity either on Σ_u and A_i ($i = 0, 1, \dots, p$) or on B_i ($i = 0, 1, \dots, p$). The matrices A_i ’s are coefficients for regressing \mathbf{Y}_t on $\mathbf{Y}_{t-1}, \dots, \mathbf{Y}_{t-p}$, while the matrices B_i ($i = 1, \dots, p$) and $I_N - B_0$ are obtained by regressing recursively an element of \mathbf{Y}_t on $\mathbf{Y}_{t-1}, \dots, \mathbf{Y}_{t-p}$ and \mathbf{Y}_t except for the own variable. In other words, the sparsity on the parameters B_0 and B_i ($i = 1, \dots, p$) corresponds to the absence of contemporaneous effects and dependence with respect to past observations, respectively. As investigated by Tunncliffe-Wilson and Reale (2008), Oxley, Reale, Tunncliffe-Wilson (2009), and Ahelegbey, Billio and Casarin (2016), Σ_u^{-1} (or B_0) and B_i ($i = 1, \dots, p$) in the SVAR model (2) can be interpreted in terms of directed acyclic graph in the graphical modeling. For this reason, we assume the sparsity on Σ_u^{-1} (or B_0) and B_i ($i = 1, \dots, p$). However, rather than fostering sparsity on the B_i ’s directly and thus apply a standard regularized estimation on a OLS/Gaussian MLE loss, we consider sparsity on the precision matrix of the vector $\mathbf{X}_t = (\mathbf{Y}_{t-p}^\top, \dots, \mathbf{Y}_{t-1}^\top, \mathbf{Y}_t^\top)^\top$, whose partial correlation coefficients characterise the coefficients among the B_i ’s and Σ_u : see, e.g., Section 5.3 of Johnston (1972). Denote $\Sigma_x = \text{Var}(\mathbf{X}_t)$, where \mathbf{X}_t is the $N(p+1)$ vector defined above. Then, denoting $\Gamma_i = \mathbb{E}[\mathbf{Y}_t \mathbf{Y}_{t-i}']$ the autocovariance matrix, we obtain

$$\Sigma_x = \begin{pmatrix} \Gamma_0 & \Gamma_1' & \cdots & \Gamma_p' \\ \Gamma_1 & \Gamma_0 & \cdots & \Gamma_{p-1}' \\ \vdots & \vdots & \ddots & \vdots \\ \Gamma_p & \Gamma_{p-1} & \cdots & \Gamma_0 \end{pmatrix}$$

Applying Example 2.2 of Dahlhaus and Eichler (2000) under the non-Gaussian assumption, equation (5) of Ahelegbey, Billio and Casarin (2016) indicates the equivalence

$$\sigma_u^{rs} = 0 \iff \text{Corr}(\mathbf{Y}_{r,t}, \mathbf{Y}_{s,t} | \{\mathbf{X}_t \setminus \{\mathbf{Y}_{r,t}, \mathbf{Y}_{s,t}\}\}) = 0,$$

$$B_{i,rs} = 0 \iff \text{Corr}(\mathbf{Y}_{r,t}, \mathbf{Y}_{s,t-i} | \{\mathbf{X}_t \setminus \{\mathbf{Y}_{r,t}, \mathbf{Y}_{s,t-i}\}\}) = 0,$$

for $i = 0, 1, \dots, p$, where σ_u^{rs} is the (r, s) th element of Σ_u^{-1} . To derive the partial correlation coefficient between two variables in \mathbf{X}_t , we need to use the inverse of Σ_x . Denoting the (j, k) th element of Σ_x^{-1} as σ_x^{jk} ,

the partial correlation coefficient between the j th and the k th elements of \mathbf{X}_t is then $\rho_x^{jk} = -\sigma_x^{jk} / \sqrt{\sigma_x^{jj} \sigma_x^{kk}}$ ($j \neq k$). It is straightforward to show

$$\sigma_u^{rs} = \sigma_x^{Np+r, Np+s}$$

for $r, s = 1, \dots, N$. By equation (5.31) of Johnston (1972), we obtain

$$B_{i,rs} = \frac{\Sigma_{x,kk}}{\Sigma_{x,jj}} \rho_x^{jk}, \quad j = Np + r, \quad k = N(p - i) + s,$$

for $i = 0, 1, \dots, p$ and $r, s = 1, \dots, N$. In view of such structure, we can impose sparsity on B_i 's and Σ_u^{-1} through the sparsity of Σ_x^{-1} , the precision matrix of the random vector \mathbf{X}_t : if we are in a position to get a consistent and a positive definite estimator of Σ_x^{-1} , say $\widehat{\Sigma}_x^{-1}$, then we would obtain consistent estimators $\widehat{B}_i, \widehat{A}_i$ ($i = 1, \dots, p$) and $\widehat{\Sigma}_u$. There is an additional merit for imposing sparsity on Σ_x^{-1} : as we consider positive definite Σ_x or Σ_x^{-1} , the parameters satisfy the stationary condition automatically by, e.g., equation (2.1.43) of Lütkepohl (2006). Such matching between the coefficients of Σ_x^{-1} and the VAR coefficients allows for a component-wise sparse structure, in the same spirit as in Ravikumar et al. (2011) or Zhang and Zou (2014). Such element-wise sparsity assumption is a key difference with the Cholesky based sparse assumption of Huang et al. (2006), which imposes a particular variable ordering as emphasised by Bickel and Levina (2008b). The sparsity assumption is thus stated for the B_i 's and Σ_u^{-1} parameters, which is equivalent to assuming a component-wise sparse structure on Σ_x^{-1} . The main contribution of our study is to propose novel sparse estimators for $\widehat{\Theta} := \widehat{\Sigma}_x^{-1}$ based on Bregman divergences and provide error bounds together with the conditions for support recovery. In particular, the von Neuman and least squares based estimators of such precision matrix have not been studied. In the rest of the paper, we denote $\Theta_0 := \Sigma_x^{-1}$ the true precision matrix of the $N(p + 1)$ random vector \mathbf{X}_t .

It is worth mentioning the competing estimation techniques for SVAR models. The inference for the 'A-model' can be performed by maximum likelihood under suitable restrictions as described in Subsection 9.3.1 of Lütkepohl (2005). Lütkepohl (2017) considered the method of moments for estimating B_0 using the OLS residuals from the VAR model. Should one consider a graphical approach for SVAR identification,

a Bayesian inference procedure can be carried out by Markov chain Monte Carlo method as in Ahelegbey, Billio and Casarin (2016). Unlike the latter work, the matching we highlighted between the SVAR parameters and the precision matrix enables us to carry out a novel sparse inference.

3 SVAR inference through sparse precision matrix

3.1 Framework

Under the sparsity assumption of $B_i, i = 0, \dots, p$ and Σ_u and thus of Σ_x^{-1} , we aim at recovering the true sparse support \mathcal{A} of the inverse of the variance covariance matrix parameter $\Theta_0 = \Sigma_x^{-1} \in \mathcal{M}_{d \times d}(\mathbb{R})$ with $d = N(p + 1)$. The sparsity assumption is specified as follows.

Assumption 1. *The true parameter $\theta_0 = \text{vec}(\Theta_0)$ is sparse so that $k_0 = \text{card}(\mathcal{A})$, where $\mathcal{A} = \{1 \leq i \leq d^2 : \theta_{i,0} \neq 0\}$ with $k_0 < d^2$ the total number of parameters.*

To estimate such sparse Θ_0 , we rely on a regularized M-estimation problem given by:

$$\hat{\Theta} = \arg \min_{\Theta \in \Omega} \left\{ \mathbb{L}_T(\Theta) + \mathbf{p}(\lambda_T, \theta) \right\}, \quad (3)$$

where $\theta = \text{vec}(\Theta)$, $\mathbf{p}(\lambda_T, \cdot) : \mathbb{R}^{d^2} \rightarrow \mathbb{R}$ is the penalty function, with λ_T the regularization parameter, which depends on the sample size, and enforce a particular type of sparse structure in the solution $\hat{\Theta}$. Here $\mathbb{L}_T : \mathbb{R}^{d^2} \times \mathbb{R}^{TN} \rightarrow \mathbb{R}$ is the non-penalized loss function, which evaluates the precision of the fit with the sample $(\mathbf{Y}_1, \dots, \mathbf{Y}_T)$. As stated in Subsections 3.2 and 5.2, the scaling behaviour of (T, d, k_0) together with statistical consistency highly depend on the choice of $\mathbb{L}_T(\cdot)$ and its regularity. Ω denotes a $d \times d$ -variance convex covariance matrix subset defined as $\Omega = \left\{ \Theta : \Theta \succ 0, \|\theta\|_1 \leq R \right\}$. The constraint through R may be somewhat arbitrary but it enforces the estimated optimum to be close enough to the theoretical optimum and ensures that $\mathbb{L}_T(\cdot)$ is lower bounded. Indeed, due to the potential non-convexity of the criterion, we include the side condition $\|\theta\|_1 \leq R$, where R is a supplementary regularization parameter to ensure the existence of local/global optima: more details on this constraint can be found in Section 3 of Loh and Wainwright (2017) or in Pognard and Fermanian (2021). Alternative presentations of the sparse

precision matrix setting are possible, where only the off-diagonal entries of Θ are penalized. Similar results for statistical consistency actually hold in this case. In our framework, we assume that all components are equally penalized to clarify our arguments. Similar settings are also considered by Loh and Wainwright (2015) and Fan, Feng and Wu (2009).

As for the loss $\mathbb{L}_T(\cdot)$ in (3), we consider the general framework of the Bregman divergence criterion $D_\phi(\Theta, \Theta_0)$, defined as a dissimilarity measure between two symmetric positive definite matrices, the candidate Θ and the true inverse Θ_0 , defined as

$$D_\phi(\Theta, \Theta_0) = \phi(\Theta) - \phi(\Theta_0) - \text{tr}((\nabla\phi(\Theta_0))^\top (\Theta - \Theta_0)),$$

where ϕ is a differentiable and strictly convex function over the space of real and symmetric positive definite matrices. Obviously, such discrepancy can not be optimized with respect to Θ unless Θ_0 is replaced by some known quantity. We propose to replace Θ_0 by its empirical version \widehat{S}^{-1} , the inverse of the sample variance covariance matrix of \mathbf{X}_t , where we assume $\widehat{S} \succ 0$. Thus, replacing Θ_0 by \widehat{S}^{-1} , we consider the loss function in (3) as $\mathbb{L}_T(\cdot) = \mathbb{L}_{T,\eta}(\cdot)$ with

$$L_{T,\eta}(\Theta) = \eta D_\phi(\Theta, S^{-1}) + (1 - \eta) D_\phi(S^{-1}, \Theta),$$

where $0 \leq \eta \leq 1$ is a known scalar value. This loss function is an extension of the standard Bregman divergence setting. Indeed, the loss function is a balance between the Bregman divergence and its switched argument version. As it will be emphasised in our theoretical analysis, suitable choices of η will enable us to provide the probability for which statistical accuracy holds. We propose the following specifications of ϕ and hence of $\mathbb{L}_T(\cdot)$:

- (i) $\phi(\Theta) = -\log(|\Theta|)$, so that the corresponding Bregman divergence can be written as

$$\mathbb{L}_{T,\eta}(\Theta) = (1 - 2\eta) \log(|\Theta|) + \text{tr}\left(\eta \widehat{S} \Theta + (1 - \eta) \Theta^{-1} \widehat{S}^{-1}\right),$$

where the terms independent of Θ are discarded. This function is known as Stein's loss, and is closely related to the standard Gaussian QML criterion up to some constants. When $\eta = 0$, one obtain the

standard Gaussian QML criterion for inverse variance covariance estimation.

(ii) $\phi(\Theta) = \text{tr}(\Theta^2)$, then the Bregman divergence is a least squares loss defined as $\mathbb{L}_{T,\eta}(\Theta) = \|\Theta - \widehat{S}^{-1}\|_F^2$, which is η independent: in that case, the dependence of $\mathbb{L}_{T,\eta}(\cdot)$ on η is simply skipped. We use $\mathbb{L}_T(\cdot)$ as the notation for the least squares case.

(iii) $\phi(\Theta) = \text{tr}(\Theta \log(\Theta) - \Theta)$, then the derivative becomes $\nabla_{\Theta} \phi(\Theta) = \text{tr}(\log(\Theta))$ (see exercise 13.31 of Abadir and Magnus, 2005, for the matrix logarithm derivative) and the Bregman divergence becomes the von Neumann. Then, our loss is defined as

$$\begin{aligned} \mathbb{L}_{T,\eta}(\Theta) = & (2\eta - 1)\text{tr}(\widehat{S}^{-1} - \Theta) + \eta\text{tr}(\Theta \log(\Theta)) + (1 - \eta)\text{tr}(\widehat{S}^{-1} \log(\widehat{S}^{-1})) \\ & - \eta\text{tr}(\log(\widehat{S}^{-1})\Theta) - (1 - \eta)\text{tr}(\log(\Theta)\widehat{S}^{-1}). \end{aligned}$$

It is interesting to compare the above losses with the D-trace loss function $\mathbb{L}_T(\Theta) = \frac{1}{2}\text{tr}(\Theta^2 \widehat{S} - 2\Theta)$ that was proposed by Zhang and Zou (2014), who performed a LASSO regularization to obtain a sparse precision matrix. We will extend their framework to non-convex penalty functions, which will in particular allow for relaxing their so-called incoherence condition when analysing the support recovery property.

To summarize, the loss function $\mathbb{L}_T(\cdot)$ in (3) will be taken as the Stein's and von Neumann losses with switched arguments $\mathbb{L}_{T,\eta}(\cdot)$, the least squares loss and the D-trace loss. As for the penalty function $\mathbf{p}(\lambda_T, \cdot)$ in (3), we rely on the following assumption.

Assumption 2. *We consider penalty functions that are assumed to be amenable regularizers defined as follows. We denote $\mathbf{p}(\cdot, \cdot) : \mathbb{R}_+ \times \mathbb{R}^q$, with q denoting the dimension problem, the penalty function - or regularizer -, which is assumed to be coordinate-separable with respect to $\theta \in \mathbb{R}^q$, idest $\mathbf{p}(\lambda_T, \theta) = \sum_{k=1}^q \mathbf{p}(\lambda_T, \theta_k)$. Furthermore, let $\mu \geq 0$, and $\mathbf{p}(\lambda_T, \cdot)$ is μ -amenable if*

(i) $x \mapsto \mathbf{p}(\lambda_T, x)$ is symmetric around zero and $\mathbf{p}(\lambda_T, 0) = 0$.

(ii) $x \mapsto \mathbf{p}(\lambda_T, x)$ is non-decreasing on \mathbb{R}^+ .

(iii) $x \mapsto \frac{\mathbf{p}(\lambda_T, x)}{x}$ is non-increasing on \mathbb{R}_*^+ .

(iv) $x \mapsto p(\lambda_T, x)$ is differentiable for any $x \neq 0$.

(v) $\lim_{x \rightarrow 0^+} \partial_x p(\lambda_T, x) = \lambda_T$.

(vi) $x \mapsto p(\lambda_T, x) + \frac{\mu}{2}x^2$ is convex for some $\mu \geq 0$.

The regularizer $\mathbf{p}(\lambda_T, \cdot)$ is (μ, ζ) -amenable if in addition

(vii) There exists $\zeta \in (0, \infty)$ such that $\partial_x p(\lambda_T, x) = 0$ for $x \geq \lambda_T \zeta$.

Let $\mathbf{q} : \mathbb{R}^+ \times \mathbb{R}^q \rightarrow \mathbb{R}$ be $\mathbf{q}(\lambda_T, x) = \lambda_T \|x\|_1 - \mathbf{p}(\lambda_T, x)$ so that the function $\frac{\mu}{2} \|x\|_2^2 - \mathbf{q}(\lambda_T, x)$ is convex.

Assumption 1 implies that the true support (unknown) is sparse, that is the matrix Θ_0 contains zero components. The regularization - or penalization - procedure provides an estimator of \mathcal{A} . To derive our theoretical properties, assumption 2 provides regularity conditions that potentially encompass non-convex functions. These regularity conditions are the same than Loh and Wainwright (2015, 2017) or Pognard and Fermanian (2021). In this paper, we focus on the LASSO, the SCAD due to Fan and Li (2001) and the MCP due to Zhang (2010), respectively defined as

$$\text{LASSO} : \mathbf{p}(\lambda_T, \rho) = \lambda_T |\rho|,$$

$$\text{MCP} : \mathbf{p}(\lambda_T, \rho) = \text{sign}(\rho) \lambda_T \int_0^{|\rho|} (1 - z/(\lambda_T b_{\text{mcp}}))_+ dz,$$

$$\text{SCAD} : \mathbf{p}(\lambda_T, \rho) = \begin{cases} \lambda_T |\rho|, & \text{for } |\rho| \leq \lambda_T, \\ -(\rho^2 - 2b_{\text{scad}} \lambda_T |\rho| + \lambda_T^2)/(2(b_{\text{scad}} - 1)), & \text{for } \lambda_T \leq |\rho| \leq b_{\text{scad}} \lambda_T, \\ (b_{\text{scad}} + 1) \lambda_T^2 / 2, & \text{for } |\rho| > b_{\text{scad}} \lambda_T, \end{cases}$$

where $b_{\text{scad}} > 2$ and $b_{\text{mcp}} > 0$ are fixed parameters for the SCAD and MCP respectively. The LASSO is a μ -amenable regularizer, whereas the SCAD and the MCP are (μ, ζ) -amenable. More precisely, $\mu = 0$ (resp. $\mu = 1/(b_{\text{scad}} - 1)$, resp. $\mu = 1/b_{\text{mcp}}$) for the Lasso (resp. SCAD, resp. MCP). The parameter μ can be interpreted as a coefficient of non-convexity level: the larger, the more non-convex the penalty becomes.

The penalized problem (3) may not be convex depending on the choice of the penalty - SCAD or MCP - and/or for a specific Bregman divergence. Therefore, we would like to weaken the convexity assumption so that we could evaluate the accuracy of $\hat{\Theta}$. To do so, the restricted strong convexity is a key ingredient

to handle non-convex loss functions. Intuitively, we would like to handle a loss function that locally admits some curvature. To ensure this property, we rely on the strong convexity (local) of the loss function. The strong convexity of a differentiable loss function corresponds to a strictly positive lower bound on the eigenvalues of the Hessian matrix uniformly valid over a local region around the true parameter. This amounts to a curvature condition. More precisely, we are interested in a particular direction, that is the difference $\Delta = \hat{\theta} - \theta_0$. Hence the notion of restricted strong convexity weakens the (local) strong convexity by adding a tolerance term. A detailed explanation is provided in Negahban et al. (2012).

Slightly extending the definition of Loh and Wainwright (2017), we say that an empirical loss function $\mathbb{G}_T(\cdot)$ satisfies the restricted strong convexity condition (RSC) at $\theta \in \mathbb{R}^q$ if there exist two positive functions α_1, α_2 and two nonnegative functions τ_1, τ_2 of (θ, T, q) such that, for any $\Delta \in \mathbb{R}^q$,

$$\begin{aligned} \langle \nabla_{\theta} \mathbb{G}_T(\theta + \Delta) - \nabla_{\theta} \mathbb{G}_T(\theta), \Delta \rangle &\geq \alpha_1 \|\Delta\|_2^2 - \tau_1 \frac{\log(q)}{T} \|\Delta\|_1^2, \text{ if } \|\Delta\|_2 \leq 1, \\ \langle \nabla_{\theta} \mathbb{G}_T(\theta + \Delta) - \nabla_{\theta} \mathbb{G}_T(\theta), \Delta \rangle &\geq \alpha_2 \|\Delta\|_2 - \tau_2 \sqrt{\frac{\log(q)}{T}} \|\Delta\|_1, \text{ if } \|\Delta\|_2 \geq 1. \end{aligned}$$

The RSC property is fundamentally local and $\alpha_k, \tau_k, k = 1, 2$ depend on the chosen θ . The RSC condition of Loh and Wainwright (2015) is similar but uniform with respect to (T, q) . Moreover, to weaken notations, we simply write α_k and $\tau_k, k = 1, 2$, by skipping their implicit arguments (θ, T, q) . The threshold for $\|\Delta\|_2$ has been set for convenience and one can reparameterize the model with $\bar{\theta} := r\theta$ for some $r > 0$.

3.2 Error bounds

We first provide some error bounds for the estimator (3) assuming that $\mathbb{L}_T(\cdot)$ satisfies the RSC condition and the penalty is μ -amenable. We assume that the population risk function $\mathbb{L}(\Theta) = \mathbb{E}[\mathbb{L}_T(\Theta)]$ is assumed to be uniquely minimized at $\theta_0 = \text{vec}(\Theta_0) \in \mathbb{R}^{d^2}$. Then we have the following Theorem.

Theorem 3.1. *Assume $\theta \in \mathbb{R}^{d^2}$ and the objective function $\mathbb{L}_T(\cdot) : \mathbb{R}^{d^2} \mapsto \mathbb{R}$ satisfies the RSC condition and $\mathbf{p}(\lambda_T, \cdot)$ is μ -amenable, with $\frac{3}{4}\mu \leq \alpha_1$. Choose*

$$4 \max \left\{ \|\nabla_{\theta} \mathbb{L}_T(\theta_0)\|_{\infty}, \alpha_2 \sqrt{\frac{\log(d^2)}{T}} \right\} \leq \lambda_T \leq \frac{\alpha_2}{6R}, \quad (4)$$

and assume $T \geq \frac{16R^2 \max\{\tau_1^2, \tau_2^2\}}{\alpha_2^2} \log(d^2)$. Let $\hat{\theta}$ be a stationary point of (3). Then $\hat{\theta}$ satisfies

$$\|\hat{\theta} - \theta_0\|_2 \leq \frac{6\lambda_T \sqrt{k_0}}{4\alpha_1 - 3\mu}, \quad \|\hat{\theta} - \theta_0\|_1 \leq \frac{6(16\alpha_1 - 9\mu)}{(4\alpha_1 - 3\mu)^2} \lambda_T k_0.$$

Remark. (i) This result is based on an optimization reasoning only and is obtained in a deterministic way; the proof can be found in Poignard and Fermanian (2021), Theorem 1. As will be clarified in the following Corollaries, to apply Theorem 3.1, we will need to check the conditions for which the loss function $\mathbb{L}_T(\cdot)$ satisfies the RSC condition. Moreover, we will show that suitable choices of λ_T and R provide the probability to satisfy the conditions of Theorem 3.1 with high probability.

(ii) About (α_1, μ) : the tightness of the error bounds are sensitive to the difference $4\alpha_1 - 3\mu$, assuming λ_T, k_0 fixed. Here, α_1 should be thought as the curvature of \mathbb{L}_T : the bigger α_1 is, the larger the curvature becomes. On the other hand, μ measures the non-convexity of the penalty function: the larger μ is, the more non-convex $\mathbf{p}(\lambda_T, \cdot)$ becomes. Thus, there is a trade-off between α_1 and μ when satisfying the constraint $4\alpha_1 > 3\mu$.

One of the purposes of the paper is to answer the following points: given $\mathbb{L}_T(\cdot)$ - Stein's and von Neumann loss with $\mathbb{L}_T := \mathbb{L}_{T,\eta}$, least squares and D-trace loss -, is the RSC condition satisfied? Can we apply Theorem 3.1 and evaluate the probability of (4)? To tackle the latter issue, we derive an exponential-type inequality. To do so, we rely on a Bernstein-type based inequality applied to the difference $\|\hat{S} - \Sigma_x\|_{\max}$. Such quantity is key when bounding the gradient of $\mathbb{L}_T(\cdot)$. The sample variance covariance matrix is defined as $\hat{S} := \frac{1}{T} \sum_{t=1}^T \mathbf{X}_t \mathbf{X}_t^\top$. Suitable assumptions on the variables are required to adapt the exponential bound to the data dependent setting. This is the motivation of the next assumptions. To do so, we consider standard assumptions on the reduced form VAR equation (1) that can be written as $\sum_{i=0}^p L(i) \mathbf{Y}_{t-i} = \mathbf{u}_t$, with $L(0) = I_N$ and $L(i) = -A_i$. We assume:

Assumption 3. *The absolute values of the zeros of the polynomial $\det(P(z))$ with $P(z) = \sum_{i=1}^p L(i) z^i, z \in \mathbb{C}$ are strictly greater than one.*

Assumption 4. *The probability distribution of \mathbf{u}_t is absolutely continuous with respect to the Lebesgue measure on \mathbb{R}^N .*

Equipped with assumptions 3 and 4, the process (\mathbf{Y}_t) is a geometrically completely regular process by Theorem 1 of Mokkadem (1988): thus (\mathbf{Y}_t) is a strongly mixing process. Hence, setting p finite, (\mathbf{X}_t) is also strongly mixing by, e.g., Theorem 3.49 of White (2001). The mixing condition will be crucial in our analysis to evaluate the probability of the discrepancy $\|\widehat{S} - \Sigma_x\|_{\max}$ exceeding a certain threshold.

Assumption 5. *(\mathbf{X}_t) is strongly mixing with mixing coefficient $\alpha(\tau) \leq \exp(-c\tau^{\gamma_1})$ with γ_1 and c positive constants. Moreover, $\exists \gamma_2 > 0, b > 0$ such that $\forall \delta > 0$ and $\forall i \leq d$, $\mathbb{P}(|\mathbf{X}_{i,t}| \geq \delta) \leq \exp(-(\frac{\delta}{b})^{\gamma_2})$.*

The latter assumption requires the exponential-type tails for the distribution of $(\mathbf{X}_{1,t}, \dots, \mathbf{X}_{d,t})$. This allows us to derive an exponential bound on $\frac{1}{T} \sum_{t=1}^T \mathbf{X}_{i,t} \mathbf{X}_{j,t} - \Sigma_{ij,x}$. We have the following lemma.

Lemma 3.2. *Let $\gamma < 1$ with $1/\gamma = 1/\gamma_1 + 3/\gamma_2$. Under assumptions 3, 4 and 5, assume $T \geq 4$, there exist positive constants C_1, C_2, C_3, C_4, C_5 depending only on b, γ_1, γ_2 such that $\forall \epsilon > 0$,*

$$\mathbb{P}(\|\widehat{S} - \Sigma_x\|_{\max} \geq \epsilon) \leq d^2 \left\{ T \exp\left(-\frac{(T\epsilon)^\gamma}{C_1}\right) + \exp\left(-\frac{(T\epsilon)^2}{C_2(1+TC_3)}\right) + \exp\left(-\frac{(T\epsilon)^2}{C_4 T} \exp\left(\frac{(T\epsilon)^{\gamma(1-\gamma)}}{C_5(\log(T\epsilon))^\gamma}\right)\right) \right\}.$$

In particular, let $L > 0$, then $\mathbb{P}(\|\widehat{S} - \Sigma_x\|_{\max} \geq L\sqrt{\frac{\log(d^2)}{T}}) = O(\frac{1}{d^2})$.

Remark. (i) This concentration inequality will be applied when bounding the random quantity $\|\mathbb{L}_{T,\eta}(\Sigma)\|_\infty$,

where the difference $\widehat{S} - \Sigma_x$ will typically be bounded in this score function.

(ii) The choice of ϵ proportional to $\sqrt{\frac{\log(d^2)}{T}}$ is motivated by condition (4) in Theorem 3.1, where we aim at evaluating the probability of satisfying such condition.

(iii) Alternatively, we can use exponential bounds for separately Lipschitz functions such as Dedecker and Fan (2015) or Alquier et al. (2020). To do so, some contraction property on the data generating process of the reduced form VAR would be necessary.

Indeed, Theorem 3.1 is stated in a deterministic manner. We show that for suitable parameter choices (λ_T, R) , the conditions of Theorem 3.1 hold with high probability. To do so, this requires bounding the random quantity $\|\nabla_{\theta} \mathbb{L}_{T,\eta}(\Theta_0)\|_{\infty}$ and verifying the RSC conditions. This motivates the use of Lemma 3.2.

We are in a position to provide the conditions for consistency of the Bregman divergence based sparse estimators. First, let us consider the Stein's loss case. For $\phi(\Theta) = -\log(|\Theta|)$, the statistical criterion is

$$\begin{cases} \hat{\Theta}^{\mathfrak{g}} &= \arg \min_{\Theta \in \Omega} \left\{ \mathbb{L}_{T,\eta}(\Theta) + \mathbf{p}(\lambda_T, \theta) \right\}, \text{ with} \\ \mathbb{L}_{T,\eta}(\Theta) &= (1 - 2\eta) \log(|\Theta|) + \text{tr}(\eta \hat{S} \Theta + (1 - \eta) \Theta^{-1} \hat{S}^{-1}). \end{cases} \quad (5)$$

Here, the loss function is optimized over the convex set $\Omega = \Omega = \left\{ \Theta : \Theta \succ 0, \|\theta\|_1 \leq R \right\}$. The side constraint in Ω is in the same spirit as in Loh and Wainwright (2015).

Corollary 3.3. *Assume the regularizer is μ -amenable, under the sample size $T \geq CR^2 \alpha_2^{-2} \log(d^2)$, with $C > 0$ a sufficiently large constant, with $\alpha_2 = (2\eta - 1) \{ \lambda_{\max}(\Theta_0) + 1 \}^{-2} + 2(1 - \eta) \{ \lambda_{\max}(\Theta_0) + 1 \}^{-3} \lambda_{\min}(\hat{S}^{-1})$, if the regularization parameter satisfies*

$$4 \max \left\{ \|(1 - 2\eta) \Theta_0^{-1} + \eta \hat{S} - (1 - \eta) \Theta_0^{-1} \hat{S}^{-1} \Theta_0^{-1}\|_{\max}, \alpha_2 \sqrt{\frac{\log(d^2)}{T}} \right\} \leq \lambda_T \leq \frac{\alpha_2}{6R}, \quad (6)$$

where $\Theta_0 \in \Omega$, suppose $\frac{3}{4}\mu < \alpha_1$ with $\alpha_1 = \alpha_2$. Then any local optimum $\hat{\Theta}^{\mathfrak{g}}$ of program (5) satisfies

$$\|\hat{\Theta}^{\mathfrak{g}} - \Theta_0\|_F \leq \frac{6\lambda_T \sqrt{k_0}}{4\alpha_1 - 3\mu}, \quad \|\text{vec}(\hat{\Theta}^{\mathfrak{g}}) - \text{vec}(\Theta_0)\|_1 \leq \frac{6(16\alpha_1 - 9\mu)\lambda_T k_0}{(4\alpha_1 - 3\mu)^2}, \quad (7)$$

with $\alpha_1 = \alpha_2$, $k_0 = |\mathcal{A}|$ and $\mathcal{A} = \{1 \leq i \leq d^2 : \theta_{i,0} := \text{vec}(\Theta_0)_i \neq 0\}$.

Furthermore, for $\eta = 1$, under assumptions 3, 4 and 5 so that the sample variance-covariance estimator satisfies the bound in Lemma 3.2, if (λ_T, R) are chosen so that $C_1 \sqrt{\log(d^2)/T} \leq \lambda_T \leq C_2/R$ and for a sample size $T \geq L \left\{ \log(d^2) \max(R^2, k_0) \vee \log(d^2)^{2/\gamma-1} \right\}$, for C_1, C_2, L large constants, then (7) hold with probability at least $1 - O(\exp(-\log(d^2))) - o(\exp(-\log(d^2)))$.

Remark. (i) The proof relies on the following two steps: first, we verify the RSC condition for $\mathbb{L}_T(\Theta)$, derive the quantities $\alpha_1, \alpha_2, \tau_1, \tau_2$ by lower bounding $\lambda_{\min}(\nabla_{\theta\theta^{\top}}^2 \mathbb{L}_T(\Theta))$ and obtain the error bounds (7) applying Theorem 3.1; second, we bound $\|\nabla_{\theta} \mathbb{L}_T(\Theta_0)\|_{\infty}$ using Lemma 3.2 for a fixed ϵ proportional to $\sqrt{\log(d^2)/T}$. In that case, the required rate becomes $\max(R^2, k_0) \log(d^2) = O(T)$.

(ii) Should we consider the Stein's loss for estimating the variance-covariance $\Sigma_x = \Theta_0^{-1}$, the regularity of the loss is significantly altered. The RSC parameters would involve the minimum eigenvalue $\lambda_{\min}(2\widehat{S} - \Sigma_x)$, which thus must be assumed positive to ensure positive RSC parameters α_1, α_2 : see, e.g., Poignard and Terada (2020), who considered the RSC property for the Gaussian QML estimator of the factor model based variance covariance matrix.

(iii) The RSC parameters are sample dependent through \widehat{S}^{-1} . Actually, using the exponential bound on $\|\widehat{S} - \Sigma_x\|_{\max}$, we could express α_1, α_2 with respect to $\lambda_{\max}(\Sigma_x) = \lambda_{\max}(\Theta_0^{-1})$ as:

$$\lambda_{\min}(\widehat{S}^{-1}) \geq \lambda_{\max}(\widehat{S})^{-1} \geq \left(\|\widehat{S} - \Sigma_x\|_s + \|\Sigma_x\|_s \right)^{-1} \geq \left(L\sqrt{d^2 \frac{\log(d^2)}{T}} + \|\Sigma_x\|_s \right)^{-1},$$

with high probability using Lemma 3.2 with $L > 0$ and for a suitable sample size.

(iv) When $\mathbf{p}(\lambda_T, \theta) = \lambda_T \|\theta\|_1$, then setting $\lambda_T \geq L\sqrt{\log(d^2)/T}$ and $R = m_0\sqrt{k_0}$ with a constant $m_0 \geq \|\theta_0\|_2$, we have the scaling $T \geq Mk_0 \log(d^2)$.

We now consider the consistency of the least squares estimator $\widehat{\Theta}^{\text{ls}}$ (case $\phi(\Theta) = \text{tr}(\Theta^2)$) that satisfies:

$$\begin{cases} \widehat{\Theta}^{\text{ls}} &= \arg \min_{\Theta \in \Omega} \left\{ \mathbb{L}_T(\Theta) + \mathbf{p}(\lambda_T, \theta) \right\}, \text{ with} \\ \mathbb{L}_T(\Theta) &= \|\Theta - \widehat{S}^{-1}\|_F^2, \quad \Omega = \left\{ \Theta : \Theta \succ 0, \|\theta\|_1 \leq R \right\}. \end{cases} \quad (8)$$

Corollary 3.4. *Assume the regularizer is μ -amenable, under the sample size $T \geq C \max(R^2, k_0^2) \alpha_2^{-2} \log(d^2)$,*

with $C > 0$ a sufficiently large constant, with $\alpha_2 = 2$, if the regularization parameter satisfies

$$4 \max \left\{ \|2(\Theta_0 - \widehat{S}^{-1})\|_{\max}, \alpha_2 \sqrt{\frac{\log(d^2)}{T}} \right\} \leq \lambda_T \leq \frac{\alpha_2}{6R}, \quad (9)$$

where $\Theta_0 \in \Omega$, suppose $\frac{3}{4}\mu < \alpha_1$ with $\alpha_1 = \alpha_2$. Then any local optimum $\widehat{\Theta}^{\text{ls}}$ of program (8) satisfies

$$\|\widehat{\Theta}^{\text{ls}} - \Theta_0\|_F \leq \frac{6\lambda_T\sqrt{k_0}}{8-3\mu}, \quad \|\text{vec}(\widehat{\Theta}^{\text{ls}}) - \text{vec}(\Theta_0)\|_1 \leq \frac{6(32-9\mu)\lambda_T k_0}{(8-3\mu)^2}, \quad (10)$$

with $k_0 = |\mathcal{A}|$ and $\mathcal{A} = \{1 \leq i \leq d^2 : \theta_{i,0} := \text{vec}(\Theta_0)_i \neq 0\}$.

Furthermore, under assumptions 3, 4 and 5 so that the sample variance-covariance estimator satisfies the bound in Lemma 3.2, if (λ_T, R) are chosen so that $C_1\sqrt{d^2 \log(d^2)/T} \leq \lambda_T \leq C_2/R$ and for a sample

size $T \geq L \left\{ d^2 \log(d^2) \max(R^2, k_0) \vee \log(d^2)^{2/\gamma-1} \right\}$, for $C_1, C_2, L > 0$ large constants, then (10) hold with probability at least $1 - O(\exp(-\log(d^2))) - o(\exp(-\log(d^2)))$.

Remark. Evaluating the probability of (9) requires upper bounding $\|\Sigma_x^{-1} - \widehat{S}^{-1}\|_{\max}$. Using the exponential bound on $\|\Sigma_x - \widehat{S}\|_{\max}$ of Lemma 3.2, one can get a bound on $\|\Sigma_x^{-1} - \widehat{S}^{-1}\|_{\max}$ at a cost d .

Let us focus on the von Neumann estimator $\widehat{\Theta}^{\text{vn}}$. To analyse its properties, we rely on the power series expansion $\log(A) = -\sum_{k=1}^{\infty} \frac{1}{k} (I_d - A)^k$, where A is a d -square symmetric positive definite matrix, to handle the matrix logarithm part. This expression is valid when $\|I_d - A\|_s < 1$. In our case, instead of assuming $\|I_d - \Theta\|_s < 1$ for expanding $\log(\Theta)$, which might be restrictive, we consider Θ/ν , where ν is the constant $\nu > \lambda_{\max}(\Theta)$ and integrate such constraint in the parameter space. Thus, we consider the criterion

$$\left\{ \begin{array}{l} \widehat{\Theta}^{\text{vn}} = \arg \min_{\Theta \in \Omega^{\text{vn}}} \left\{ \mathbb{L}_{T,\eta}(\Theta) + \mathbf{p}(\lambda_T, \theta) \right\}, \text{ with} \\ \mathbb{L}_{T,\eta}(\Theta) = (2\eta - 1) \text{tr} \left(\widehat{S}^{-1} - \Theta \right) / \nu + \eta \text{tr} \left(\Theta \log(\Theta/\nu) \right) / \nu + (1 - \eta) \text{tr} \left(\widehat{S}^{-1} \log(\widehat{S}^{-1}/\nu) \right) / \nu \\ \quad - \eta \text{tr} \left(\log(\widehat{S}^{-1}/\nu) \Theta \right) / \nu - (1 - \eta) \text{tr} \left(\log(\Theta/\nu) \widehat{S}^{-1} \right) / \nu. \end{array} \right. \quad (11)$$

where Ω^{vn} is defined as $\Omega^{\text{vn}} = \left\{ \Theta : \Theta \succ 0, \|\Theta\|_s < \nu, \|\theta\|_1 \leq R \right\}$.

Corollary 3.5. Assume the regularizer is μ -amenable, under the sample size $T \geq CR^2 \alpha_2^{-2} \log(d^2)$, with $C > 0$ large enough, let $k_0 = |\mathcal{A}|$ and $\mathcal{A} = \{1 \leq i \leq d^2 : \theta_{i,0} := \text{vec}(\Theta_0)_i \neq 0\}$, then

(i) $\eta = 1$: $\alpha_2 = 1/(\nu d)$, if the regularization parameter satisfies

$$4 \max \left\{ \left\| \frac{1}{\nu} \left(\log(\Theta_0/\nu) - \log(\widehat{S}^{-1}/\nu) \right) \right\|_{\max}, \alpha_2 \sqrt{\frac{\log(d^2)}{T}} \right\} \leq \lambda_T \leq \frac{\alpha_2}{6R}, \quad (12)$$

with $\Theta_0 \in \Omega^{\text{vn}}$, suppose $\frac{3}{4}\mu < \alpha_1$ with $\alpha_1 = \alpha_2$. Then any local optimum $\widehat{\Theta}^{\text{vn}}$ of (11) satisfies

$$\|\widehat{\Theta}^{\text{vn}} - \Theta_0\|_F \leq \frac{6\lambda_T \sqrt{k_0}}{4/(\nu d) - 3\mu}, \quad \|\text{vec}(\widehat{\Theta}^{\text{vn}}) - \text{vec}(\Theta_0)\|_1 \leq \frac{6(16/(\nu d) - 9\mu)\lambda_T k_0}{(4/(\nu d) - 3\mu)^2}. \quad (13)$$

(ii) $\eta = 0$: $\alpha_2 = \lambda_{\min}(\widehat{S}^{-1}/\nu) \{ \lambda_{\max}(\Theta_0) + 1 \}^{-2}$, if the regularization parameter satisfies

$$4 \max \left\{ \left\| \frac{1}{\nu} \Theta_0^{-1} \left(\Theta_0 - \widehat{S}^{-1} \right) \right\|_{\max}, \alpha_2 \sqrt{\frac{\log(d^2)}{T}} \right\} \leq \lambda_T \leq \frac{\alpha_2}{6R}, \quad (14)$$

with $\Theta_0 \in \Omega$, suppose $\frac{3}{4}\mu < \alpha_1$ with $\alpha_1 = \alpha_2$. Then any local optimum $\widehat{\Theta}^{\text{vn}}$ of (11) satisfies

$$\begin{aligned} \|\widehat{\Theta}^{\text{vn}} - \Theta_0\|_F &\leq \frac{6\lambda_T\sqrt{k_0}}{4\left(\lambda_{\min}(\widehat{S}^{-1}/\nu)\{\lambda_{\max}(\Theta_0) + 1\}^{-2}\right) - 3\mu}, \\ \|\text{vec}(\widehat{\Theta}^{\text{vn}}) - \text{vec}(\Theta_0)\|_1 &\leq \frac{6(16\left(\lambda_{\min}(\widehat{S}^{-1}/\nu)\{\lambda_{\max}(\Theta_0) + 1\}^{-2}\right) - 9\mu)\lambda_T k_0}{(4\left(\lambda_{\min}(\widehat{S}^{-1}/\nu)\{\lambda_{\max}(\Theta_0) + 1\}^{-2}\right) - 3\mu)^2}. \end{aligned} \quad (15)$$

Furthermore, under assumptions 3, 4 and 5 so that the sample variance-covariance estimator satisfies the bound in Lemma 3.2, if (λ_T, R) are chosen so that $C_1\sqrt{d^2 \log(d^2)/T} \leq \lambda_T \leq C_2/R$ and for a sample size $T \geq L\left\{\|\Theta_0^{-1}/\nu\|_s^2 d^2 \log(d^2) \max(R^2, k_0) \vee \log(d^2)^{2/\gamma-1}\right\}$, for $C_1, C_2, L > 0$ large constants, then (15) hold with probability at least $1 - O(\exp(-\log(d^2))) - o(\exp(-\log(d^2)))$.

Remark. (i) As in Corollary 3.3, the RSC parameters are sample dependent through \widehat{S}^{-1} and can be expressed with respect to $\lambda_{\max}(\Sigma_x)$.

(ii) For $\eta = 1$, $\nabla_{\Theta} \mathbb{L}_{T,1}(\Theta_0) = \frac{1}{\nu} \left(\log(\Theta_0/\nu) - \log(\widehat{S}^{-1}/\nu) \right)$. Then upper bounding $\log(\Theta_0/\nu) - \log(\widehat{S}^{-1}/\nu)$ is challenging. For instance, should we use power series expansion, we would need to control for

$$\left\| \sum_{k=1}^{\infty} \frac{1}{k} \left\{ (I_d - \widehat{S}^{-1}/\nu)^k - (I_d - \Theta_0/\nu)^k \right\} \right\|_{\max} \leq \sum_{k=1}^{\infty} \frac{1}{k} \left\| (I_d - \widehat{S}^{-1}/\nu)^k - (I_d - \Theta_0/\nu)^k \right\|_s.$$

Thus, unless we assume further conditions such as Θ_0 and \widehat{S}^{-1} commuting, deriving an exponential bound is a challenging task.

(iii) When $\eta = 0$, the scaling behaviour (T, d, k_0) in the von Neumann case is similar to the least squares case. Indeed, evaluating (14) requires controlling for $\frac{1}{\nu}(\Theta_0 - \widehat{S}^{-1}) = \frac{1}{\nu}(\Sigma_x^{-1} - \widehat{S}^{-1})$. For this reason, we are able to provide the probability for which inequalities in (15) hold.

We finally investigate the D-trace loss based sparse estimator $\widehat{\Theta}^{\text{dt}}$, which satisfies the criterion

$$\begin{cases} \widehat{\Theta}^{\text{dt}} &= \arg \min_{\Theta \in \Omega} \left\{ \mathbb{L}_T(\Theta) + \mathbf{p}(\lambda_T, \theta) \right\}, \\ \mathbb{L}_T(\Theta) &= \frac{1}{2} \text{tr}(\Theta^2 \widehat{S}) - \text{tr}(\Theta), \quad \Omega = \Omega = \left\{ \Theta : \Theta \succ 0, \|\theta\|_1 \leq R \right\}. \end{cases} \quad (16)$$

This loss function was proposed by Zhang and Zou (2014). We then have the following error bounds.

Corollary 3.6. Assume the regularizer is μ -amenable, under the sample size $T \geq C \max(R^2, k_0^2) \alpha_2^{-2} \log(d^2)$,

with $C > 0$ a sufficiently large constant, with $\alpha_2 = \lambda_{\min}(\widehat{S})$, if the regularization parameter satisfies

$$4 \max \left\{ \frac{1}{2} \Theta_0 (\widehat{S} - \Theta_0^{-1}) + \frac{1}{2} (\widehat{S} - \Theta_0^{-1}) \Theta_0 \right\}_{\max, \alpha_2} \sqrt{\frac{\log(d^2)}{T}} \leq \lambda_T \leq \frac{\alpha_2}{6R}, \quad (17)$$

where $\Theta_0 \in \Omega$, suppose $\frac{3}{4}\mu < \alpha_1$ with $\alpha_1 = \alpha_2$, then any local optimum $\widehat{\Theta}^{\text{dt}}$ of program (16) satisfies

$$\|\widehat{\Theta}^{\text{dt}} - \Theta_0\|_F \leq \frac{6\lambda_T \sqrt{k_0}}{4\lambda_{\min}(\widehat{S}) - 3\mu}, \quad \|\text{vec}(\widehat{\Theta}^{\text{dt}}) - \text{vec}(\Theta_0)\|_1 \leq \frac{6(16\lambda_{\min}(\widehat{S}) - 9\mu)\lambda_T k_0}{(4\lambda_{\min}(\widehat{S}) - 3\mu)^2}, \quad (18)$$

with $k_0 = |\mathcal{A}|$ and $\mathcal{A} = \{1 \leq i \leq d^2 : \theta_{i,0} := \text{vec}(\Theta_0)_i \neq 0\}$.

Furthermore, under assumptions 3, 4 and 5 so that the sample variance-covariance estimator satisfies the bound in Lemma 3.2, if (λ_T, R) are chosen so that $C_1 \sqrt{d^2 \log(d^2)/T} \leq \lambda_T \leq C_2/R$ and for a sample size $T \geq L \left\{ \|\Theta_0\|_S^2 d^2 \log(d^2) \max(R^2, k_0) \vee \log(d^2)^{2/\gamma-1} \right\}$, for $C_1, C_2, L > 0$ large constants, then (18) hold with probability at least $1 - O(\exp(-\log(d^2))) - o(\exp(-\log(d^2)))$.

Remark. (i) As in remark (i) following Corollary 3.5, the RSC parameters can be expressed with respect to $\lambda_{\min}(\Sigma_x)$ using the exponential bound of Lemma 3.2.

(ii) Zhang and Zou (2014), in their Theorem 2, obtained a similar rate for the $\|\cdot\|_F$ consistency of the D-trace estimator; and our scaling behaviour (T, d, k_0) shares the same sample size order as theirs. However, contrary to their result, our upper bound explicitly links the regularity of the D-trace loss with the non-convexity degree of the penalty function through the RSC property. In another context, using the D-Trace loss function, Wu and Li (2020) focused on the consistency of the LASSO penalized precision matrix difference $\widehat{\Delta} = \widehat{\Theta}_Y - \widehat{\Theta}_X$, where $\Theta_Y = \Sigma_Y^{-1}$, $\Theta_X = \Sigma_X^{-1}$ are the precision matrices of the two state vectors Y, X assumed sub-Gaussian. In their Theorem 2, they also obtained a similar $\|\cdot\|_F$ consistency but with much larger dimension dependent constants.

All else being equal, the curvature of the loss $\mathbb{L}_T(\cdot)$ in (3) significantly influences how informative the theoretical upper bounds are. For the Stein's loss, the curvature parameter α_1 involves the spectral norm of the true parameter Θ_0 : $\lambda_{\max}(\Theta_0)$ determines the convexity degree of $\mathbb{L}_T(\cdot)$. As it will be highlighted in

Section 4, a small α_1 value requires a low degree of non-convexity in the penalty to satisfy $4\alpha_1 > 3\mu$: the smaller μ is, the less non-convex $\mathbf{p}(\lambda_T, \cdot)$ should be. The dimension impacts also the von Neumann and D-trace losses. Interestingly, the least squares loss involves a dimension free expression for α_1 .

3.3 Support recovery

Based on the Karush-Kuhn-Tucker optimality conditions, Wainwright (2009) developed the primal dual witness (PDW) approach to derive selection consistency for convex problems. There exist similar approaches in Candès and Plan (2009) or Zhao and Yu (2006). The PDW approach consists in plugging the true subset model \mathcal{A} in the KKT optimality conditions, which are necessary and sufficient if the problem is convex, and checking if they can be satisfied. It means that any solution of the non restricted problem (the original problem providing \mathcal{A}) is also a solution to the restricted problem (the regularized one). Loh and Wainwright (2017) showed that this approach can be extended to a nonconvex problem and thus to any stationary point, which is their key contribution. They prove that all stationary points are consistent for variable selection via a strict dual feasibility condition and second-order conditions. To obtain the support recovery property, the RSC condition of the loss function with parameters $(\alpha_k, \tau_k)_{k=1,2}$ and the μ -amenability of the penalty are key assumptions. More details can be found in Subsection 5.2: there, in Theorem 5.1, we provide the conditions of Loh and Wainwright (2017) to ensure the success of the PDW construction - corresponding to **Step 3.** -, that is the scaling of (λ_T, R) and the so-called strict feasibility condition, which characterize the solution of the PDW construction; Theorem 5.2 establishes the support recovery property together with consistency in the $\|\cdot\|_\infty$ -sense under the RSC condition, μ -amenable penalties and strict dual feasibility; finally, two sufficient conditions in Proposition 5.3 ensure that strict dual feasibility holds for (μ, ζ) -amenable penalties. Within this setting, we provide ℓ_∞ -guarantees for the regularized $\hat{\Theta}$ together with the conditions to satisfy the support recovery property: for the Stein's loss, we restrict our analysis to $\eta = 1$, the most commonly used loss for sparse precision matrix estimation; $\eta = 0$ for the von Neumann case. Rather than stating the support recovery property in a deterministic manner, we directly evaluate

the probability of satisfying the latter property. For all Bregman losses, we show that any local/global optimum of (3) corresponds to the oracle estimator with high probability. The latter is given as

$$\widehat{\Theta}^{\mathcal{O}} := \arg \min_{\theta \in \mathbb{R}^{|\mathcal{A}|}: \Theta \in \Omega} \left\{ \mathbb{L}_T(\Theta) \right\}, \quad (19)$$

with $\text{vec}(\widehat{\Theta}^{\mathcal{O}}) = (\text{vec}(\widehat{\Theta}_{\mathcal{A}}^{\mathcal{O}}), \mathbf{0}_{\mathcal{A}^c})$. We denote the Fisher information matrix $\mathbf{K}_0 = \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]$. The conditions we derive hold for all stationary points of (3), idest for local/global optimum. $\mathbb{L}_T(\Theta)$ is the non-penalized loss of problem (3). Thus, we denote $\widehat{\Theta}^{\text{g},\mathcal{O}}, \widehat{\Theta}^{\text{ls},\mathcal{O}}, \widehat{\Theta}^{\text{vn},\mathcal{O}}, \widehat{\Theta}^{\text{dt},\mathcal{O}}$ the oracle estimators of (19) respectively for the Stein, least squares, von Neumann and D-trace loss.

For the Stein's loss, we consider $\eta = 1$ so that the Fisher information matrix becomes $\mathbf{K}_0 = \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)] = (\Theta_0^{-1} \otimes \Theta_0^{-1})$. The conditions for support recovery for the Stein case are given as follows:

Corollary 3.7. *Assume $T \geq L \left\{ \|\Sigma_x\|_\infty^4 k_0^2 \log(d^2) \vee \log(d^2)^{2/\gamma-1} \right\}$ with $L > 0$ large enough, the regularization parameters (λ_T, R) are chosen so that $\|\text{vec}(\Theta_0)\|_1 \leq \frac{R}{2}$ and $C_1 \sqrt{\frac{\log(d^2)}{T}} \leq \lambda_T \leq \frac{C_2}{R}$, for $C_1, C_2 > 0$, assume $\|\mathbf{K}_0^{-1}\|_\infty \leq \beta_\infty$ and assumptions 3, 4 and 5 hold. Then:*

(i) *Assume $\mathbf{p}(\lambda_T, \cdot)$ is μ -amenable penalty and $\|(\Theta_0^{-1} \otimes \Theta_0^{-1})_{\mathcal{A}^c \mathcal{A}} (\Theta_0^{-1} \otimes \Theta_0^{-1})_{\mathcal{A} \mathcal{A}}^{-1}\|_\infty \leq \omega < 1$ (incoherence condition), then with probability at least $1 - O(\exp(-\log(d^2)) - o(\exp(-\log(d^2))))$, the objective function (5) admits a unique optimum so that $\widehat{\mathcal{A}} \subseteq \mathcal{A}$ and for a sufficiently large $\tilde{L} > 0$*

$$\|\widehat{\Theta}^{\text{g}} - \Theta_0\|_{\max} \leq \tilde{L} \sqrt{\frac{\log(d^2)}{T}} + \lambda_T \beta_\infty.$$

(ii) *Assume $\mathbf{p}(\lambda_T, \cdot)$ is (μ, ζ) -amenable and $\min_{i \in \mathcal{A}} |\text{vech}(\Theta_0)_i| \geq \lambda_T(\zeta + 2\beta_\infty) + \tilde{L} \sqrt{\frac{\log(d^2)}{T}}$ for a sufficiently large $\tilde{L} > 0$, then with probability $1 - O(\exp(-\log(d^2)) - o(\exp(-\log(d^2))))$, (5) admits a unique optimum $\widehat{\Theta}^{\text{g}}$, which agrees with the oracle estimator $\widehat{\Theta}^{\text{g},\mathcal{O}}$ so that*

$$\|\widehat{\Theta}^{\text{g}} - \Theta_0\|_{\max} \leq \tilde{L} \sqrt{\frac{\log(d^2)}{T}}.$$

Remark. (i) The proof relies on the use of Theorem 5.2. To do so, strict dual feasibility must be proved (since Theorem 5.2 relies on the conditions of Theorem 5.1 and strict dual feasibility). To establish

strict dual feasibility, we use Theorem 5.1 for μ -amenable penalty functions and Proposition 5.3 for (μ, ζ) -amenable penalty functions. The main proof steps can be summarized as:

- (a) Establishing strict dual feasibility by upper bounding the quantities $\|\nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}}\|_{\infty}$ and $\|\hat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \hat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)\|_{\infty}$ by $(1 - \delta)/2\lambda_T$ for $\delta \in [0, 1]$ defined in Theorem 5.1 - with $\tau_1 = 0$ since the RSC condition for the Gaussian loss is satisfied with $\tau_1 = \tau_2 = 0$; here $\hat{\mathbf{K}}$ is defined as in Theorem 5.2 for the Gaussian loss. These bounds correspond to inequalities (27) and (28) in Proposition 5.3. Note that for the μ -amenable penalty case, the additional quantity $\|\hat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \hat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1}\|_{\infty}$ must be upper-bounded. Once strict dual feasibility is established, we compute the upper bound of $\|\hat{\Theta}^g - \Theta_0\|_{\max}$ in point (i) of Theorem 5.2.
- (b) Establishing point (ii) of Corollary 3.7 uses the exact same steps as in point (i), except that the (μ, ζ) -amenability allows for a simplification in the upper bound of $\|\hat{\Theta}^g - \Theta_0\|_{\max}$: the term involving $\|\hat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \hat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1}\|_{\infty}$ can be discarded as well as the incoherence condition.
- (ii) Corollary 3.7 is not expressed in a deterministic manner since we evaluate the probability of satisfying the inequalities (24) and (25) in Theorem 5.1. This implies controlling, among others, for the infimum norm of the score function evaluated at Θ_0 . We thus obtain λ_T proportional to $\sqrt{\log(d^2)/T}$.
- (iii) The scaling behaviour we obtained is in the same vein as in Loh and Wainwright (2017) or Ravikumar et al. (2011). Loh and Wainwright (2017) require $n > L \|\Gamma_{\mathcal{A} \mathcal{A}}^{-1}\|_{\infty}^2 \|\Sigma_x\|_{\infty}^6 r_0^2 \log(d^2)$ in their Corollary 4 for i.i.d. data, for $L > 0$ a large constant, where $\|\cdot\|_{\infty}$ denotes the ℓ_{∞} operator norm, r_0 the number of nonzero entries per row and $\Gamma = \nabla_{\theta\theta}^2 \mathbb{L}_{T,1}(\Theta_0)$. Such scaling is obtained for the side constraint $\|\Theta\|_s \leq \kappa$ rather than $\|\text{vec}(\Theta)\|_1 \leq R$, which allows their sparsity assumption to be stated at a row/column level rather than over the whole matrix, which is the main difference with our setting.
- (iv) We emphasize that there is an alternative method for constructing $\text{vec}(\hat{\Theta}^g)_{\mathcal{A}}$ such that $\text{supp}(\hat{\Theta}^g)_{\mathcal{A}} \subseteq \mathcal{A}$ and $\text{vec}(\hat{\Theta}^g)_{\mathcal{A}}$ is a zero-subgradient point of the program (22) in **Step 1** of the Primal Dual Witness method: this method is based on the Brouwer's fixed point Theorem. Intuitively, the idea is to prove

that if there is a zero sub-gradient vector of the penalized estimator $\mathbb{L}_{T,1}(\Theta) + \mathbf{p}(\lambda_T; \text{vec}(\Theta))$ within the set $\{\Theta \in \Omega^g, \text{supp}(\Theta) \subseteq \text{supp}(\Theta_0)\}$, then this vector is the unique optimum. Then Brouwer's fixed point Theorem is used to show that such optimum lies in a neighbourhood of the true value $\text{vech}(\Theta_0)$ in the $\|\cdot\|_\infty$ -sense. Such method was developed by Ravikumar et al. (2011) for LASSO penalized Θ or Loh and Wainwright (2017) for LASSO/SCAD/MCP penalized Θ .

We now consider the least squares estimator $\widehat{\Theta}^{\text{ls}}$. The oracle $\widehat{\Theta}^{\text{ls}, \mathcal{O}}$ satisfies (19) with $\mathbb{L}_T(\Theta) = \|\Theta - \widehat{S}^{-1}\|_F^2$.

Corollary 3.8. *Under assumptions 3, 4 and 5, assume $T \geq L \max\{[(d^2 - k_0) \vee k_0] \log(d^2), \log(d^2)^{2/\gamma-1}\}$ with $L > 0$ large enough, assume the regularization parameters (λ_T, R) are chosen so that $\|\text{vec}(\Theta_0)\|_1 \leq \frac{R}{2}$ and $C_1 \sqrt{\frac{\log(d^2)}{T}} \leq \lambda_T \leq \frac{C_2}{R}$ with $C_1, C_2 > 0$. Assume $\mathbf{p}(\lambda_T, \cdot)$ is a (μ, ζ) -amenable penalty and $\min_{i \in \mathcal{A}} |\text{vech}(\Theta_0)_i| \geq \lambda_T(\zeta + 2) + L \sqrt{\frac{\log(d^2)}{T}}$ for $L > 0$ large enough, then with probability $1 - O(\exp(-\log(d^2))) - o(\exp(-\log(d^2)))$, (8) admits a unique optimum $\widehat{\Theta}^{\text{ls}}$, which agrees with the oracle estimator $\widehat{\Theta}^{\text{ls}, \mathcal{O}}$ so that*

$$\|\widehat{\Theta}^{\text{ls}} - \Theta_0\|_{\max} \leq L \sqrt{\frac{\log(d^2)}{T}}.$$

Remark. (i) The proof follows the same steps as in Corollary 3.7: establishing inequalities (27) and (28) in Proposition 5.3 to use Theorem 5.2. Note that upper bounding $\|\nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}^c}\|_\infty$ and $\|\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}}\|_\infty$ by $(1 - \delta)/2\lambda_T$, for $\delta \in [0, 1]$ defined in Theorem 5.1 - with $\tau_1 = 0$ since the RSC condition for the least squares loss is satisfied with $\tau_1 = \tau_2 = 0$, is more straightforward due to the linearity of the least squares loss, contrary to the case $\phi(\Theta) = -\log(|\Theta|)$.

(ii) If the Lasso is considered, which is a μ -amenable regularizer, then the mutual incoherence condition does not hold since $\|\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1}\|_\infty = 1$. Consequently, strict dual feasibility can not be established for μ -amenable penalties when the least squares type loss function is considered.

Let us consider the conditions for support recovery for the von Neumann case, where we restrict our analysis to $\eta = 0$ to clarify our arguments. The matrix $\mathbf{K}_0 = \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]$ is given by (31) in the Appendix.

Corollary 3.9. Assume $T \geq L \max \left\{ \frac{1}{\nu^2} \left[(d^2 - k_0) \vee \frac{k_0^9}{\nu^2} \|\Sigma_x\|_s^4 \left(1 - \|\Theta_0/\nu\|_s \right)^{-2} \right] \|\Sigma_x\|_s^2 \log(d^2), \log(d^2)^{2/\gamma-1} \right\}$

with $L > 0$ large enough, assume $\mathbb{E}[\widehat{S}^{-1}] < \infty$, choose the regularization parameters (λ_T, R) as $\|\text{vec}(\Theta_0)\|_1 \leq \frac{R}{2}$ and $C_1 \sqrt{\frac{\log(d^2)}{T}} \leq \lambda_T \leq \frac{C_2}{R}$, suppose $\|\mathbf{K}_0^{-1}\|_\infty \leq \beta_\infty$ and assumptions 3, 4 and 5 hold. Then:

- (i) Assume $\mathbf{p}(\lambda_T, \cdot)$ is a μ -amenable penalty and $\|\mathbf{K}_{0, \mathcal{A}^c \mathcal{A}} \mathbf{K}_{0, \mathcal{A} \mathcal{A}}^{-1}\|_\infty \leq \omega < 1$ (incoherence condition), then with probability at least $1 - O(\exp(-\log(d^2))) - o(\exp(-\log(d^2)))$, the objective function (11) admits a unique optimum so that $\widehat{\mathcal{A}} \subseteq \mathcal{A}$ and for a sufficiently large $\tilde{L} > 0$

$$\|\widehat{\Theta}^{\text{vn}} - \Theta_0\|_{\max} \leq \tilde{L} \sqrt{\frac{\log(d^2)}{T}} + \lambda_T \beta_\infty.$$

- (ii) Assume $\mathbf{p}(\lambda_T, \cdot)$ is (μ, ζ) -amenable and $\min_{i \in \mathcal{A}} |\text{vech}(\Theta_0)_i| \geq \lambda_T (\zeta + 2\beta_\infty) + \tilde{L} \sqrt{d^2 \frac{\log(d^2)}{T}}$ for a sufficiently large $\tilde{L} > 0$, then with probability $1 - O(\exp(-\log(d^2))) - o(\exp(-\log(d^2)))$, (11) admits a unique optimum $\widehat{\Theta}^{\text{vn}}$, which agrees with the oracle estimator $\widehat{\Theta}^{\text{vn}, \mathcal{O}}$ so that

$$\|\widehat{\Theta}^{\text{vn}} - \Theta_0\|_{\max} \leq \tilde{L} \sqrt{\frac{\log(d^2)}{T}}.$$

Remark. The scaling behaviour (T, d, k_0) is less favorable compared to the Stein's loss: this is because the Hessian of the Von Neumann loss requires the control of $\|\widehat{S}^{-1}/\nu\|$, a quantity, which does not appear in the Hessian of the Stein's loss for $\eta = 1$.

Finally, we consider the D-trace loss case. The Fisher information matrix is $\mathbf{K}_0 = \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)] = \frac{1}{2} \left(\Theta_0^{-1} \otimes I_d + I_d \otimes \Theta_0^{-1} \right)$. Our analysis extends the setting of Zhang and Zou (2014) to non-convex penalties allowing for relaxing their incoherence condition in point (ii) of the following Corollary.

Corollary 3.10. Assume $T > L \max \left\{ k_0^2 \left[\|\Sigma_x\|_\infty^2 \vee \|\Theta_0\|_s^2 \right] \log(d^2), (d^2 - k_0) \|\Theta_0\|_s^2 \log(d^2), \log(d^2)^{2/\gamma-1} \right\}$,

choose the regularization parameters (λ_T, R) so that $\|\text{vec}(\Theta_0)\|_1 \leq \frac{R}{2}$ and $C_1 \sqrt{\frac{\log(d^2)}{T}} \leq \lambda_T \leq \frac{C_2}{R}$ for $C_1, C_2 > 0$, assume $\|\mathbf{K}_0^{-1}\|_\infty \leq \beta_\infty$ and assumptions 3, 4 and 5 hold. Then:

- (i) Assume $\mathbf{p}(\lambda_T, \cdot)$ is μ -amenable and $\left\| \left(\Theta_0^{-1} \otimes I_d + I_d \otimes \Theta_0^{-1} \right)_{\mathcal{A}^c \mathcal{A}} \left(\Theta_0^{-1} \otimes I_d + I_d \otimes \Theta_0^{-1} \right)_{\mathcal{A} \mathcal{A}}^{-1} \right\|_\infty \leq \omega < 1$ (incoherence condition), then with probability at least $1 - O(\exp(-\log(d^2))) - o(\exp(-\log(d^2)))$, the

objective function (16) admits a unique optimum with $\widehat{\mathcal{A}} \subseteq \mathcal{A}$ and for $L > 0$ large enough

$$\|\widehat{\Theta}^{\text{dt}} - \Theta_0\|_{\max} \leq L\sqrt{\frac{\log(d^2)}{T}} + \lambda_T\beta_\infty.$$

(ii) Assume $\mathbf{p}(\lambda_T, \cdot)$ is (μ, ζ) -amenable and $\min_{i \in \mathcal{A}} |\text{vech}(\Theta_0)_i| \geq \lambda_T(\zeta + 2\beta_\infty) + L\sqrt{\frac{\log(d^2)}{T}}$ for a sufficiently large $\tilde{L} > 0$, then with probability $1 - O(\exp(-\log(d^2))) - o(\exp(-\log(d^2)))$, (16) admits a unique optimum $\widehat{\Theta}^{\text{dt}}$, which agrees with the oracle estimator $\widehat{\Theta}^{\text{dt}, \mathcal{O}}$ so that

$$\|\widehat{\Theta}^{\text{dt}} - \Theta_0\|_{\max} \leq \tilde{L}\sqrt{\frac{\log(d^2)}{T}}.$$

Remark. Our scaling involves k_0 , which represents the total sparsity over the whole precision matrix. In their Theorem 2, Zhang and Zhou (2014) obtain a similar $\|\cdot\|_{\max}$ consistency rate, but their scaling condition for support recovery for i.i.d. data is such that $n > Lr_0^2 \log(d)$, with r_0 the maximum number of nonzero off-diagonal entries in any row (or column), a rate similar to Corollary 4 of Loh and Wainwright (2017).

4 Empirical applications

In all our simulation experiments, we simulate the N -dimensional VAR model

$$\forall 1 \leq t \leq T, \mathbf{Y}_t = \sum_{i=1}^p A_i \mathbf{Y}_{t-i} + \mathbf{u}_t, \quad \mathbf{u}_t \sim \mathcal{N}_{\mathbb{R}^N}(0, \Sigma_u), \quad (20)$$

where the matrix coefficients A_i 's are deduced from the sparse matrices B_0 and $B_i, 1 \leq i \leq p$. The matrix coefficients are simulated so that the conditions of assumption 3 are satisfied. The object of interest is the sparse inverse variance covariance matrix Θ_0 , deduced from the sparse Σ_u^{-1}, B_i 's coefficients. To recover such sparse element, we consider problems (5), (8), (11) and (16) with SCAD, MCP and LASSO penalization, providing 12 estimators. To solve the penalized optimization problem under the positive-definite constraint, we apply the numerical optimization *fmincon* on Matlab with \widehat{S}^{-1} for initialization. Alternatively, one can follow the composite gradient descent procedure of Loh and Wainwright (2015) - see their section 4 -, which consists in a three step updating procedure of the optimized parameter. To manage

the positive-definiteness constraint of Θ , one could include an alternating direction method of multipliers step in the same spirit than the ADMM algorithm provided in Appendix 3 of Bien and Tibshirani (2011).

4.1 Sensitivity analysis

We propose a sensitivity analysis of the statistical consistency and the theoretical error bound with respect to λ_T . We consider two penalization rates: $4.5\sqrt{\frac{\log(d^2)}{T}}$ and $6.5\sqrt{\frac{\log(d^2)}{T}}$, where $d = N(p+1)$. As for the side constraint parameter R , we select $R = \frac{2}{\lambda_T}\mathbf{p}(\lambda_T, \text{vec}(\Theta_0))$ to ensure the feasibility of Θ_0 following Loh and Wainwright (2015, 2017). For general data sets, R cannot be computed since the true underlying model is unknown, so that a data-driven method such as cross-validation is required. We set $N = 15$ and $p = 2$ so that $d = 45$ and the total number of parameters in Θ is 2025 (idest 1035 distinct elements). The total number of zero components is 1508 so that $k_0 = 517$, and $\|\text{vec}(\Theta_0)\|_1 = 91.12$, $\|\Theta_0\|_F = 7.81$, $\|\Theta_0\|_{\max} = 1.013$, $\|\Theta_0\|_s = 2.101$. The Stein's (resp. von Neumann) loss is calibrated for $\eta = 1$ (resp. $\eta = 0$). As for the RSC parameters, to satisfy $4\alpha_1 > 3\mu$, we consider the setting:

- (i) *Stein's loss*: $b_{\text{scad}} = 30, b_{\text{mcp}} = 20$. When $\eta = 1$, then $\alpha_1 = (\lambda_{\max}(\Theta_0) + 1)^{-2} = 0.104$. For the SCAD, $4\alpha_1 - 3\mu = 0.313$; for the MCP, $4\alpha_1 - 3\mu = 0.266$; for the Lasso, $4\alpha_1 - 3\mu = 0.4159$.
- (ii) *Least squares loss*: $\alpha_1 = 2$. We choose $b_{\text{scad}} = 30, b_{\text{mcp}} = 20$. For the SCAD, $4\alpha_1 - 3\mu = 7.896$; for the MCP $4\alpha_1 - 3\mu = 7.850$; for the Lasso, $4\alpha_1 - 3\mu = 8$.
- (iii) *von Neumann loss*: for $\eta = 0$, then $\alpha_1 = \lambda_{\min}(\widehat{S}^{-1}/\nu)(\lambda_{\max}(\Theta_0) + 1)^{-2} = 0.0167$, with $d = 45$. We set $\nu = 1$ for the sake of clarification to compute the theoretical bounds. To evaluate $\lambda_{\min}(\widehat{S}/\nu)$, we simulated 100 times the DGP for $T = 20000$ and took the average of these hundred minimum eigenvalues, so that $\lambda_{\min}(\widehat{S}^{-1}/\nu) = 0.079$ and $\alpha = 0.0082$. To satisfy $4\alpha_1 > 3\mu$, we chose $b_{\text{scad}} = 110, b_{\text{mcp}} = 100$ so that: for the SCAD, $4\alpha_1 - 3\mu = 0.0051$; for the MCP $4\alpha_1 - 3\mu = 0.0026$; for the Lasso, $4\alpha_1 - 3\mu = 0.0326$.
- (iv) *D-trace loss*: $\alpha_1 = \lambda_{\min}(\widehat{S})$; based on 100 simulations of the DGP for $T = 20000$, we obtained 100 estimates \widehat{S} and computed the average of these 100 minimum eigenvalues and obtained $\alpha_1 = 0.557$.

We set $b_{\text{scad}} = 30, b_{\text{mcp}} = 20$. For the SCAD, $4\alpha_1 - 3\mu = 2.005$; for the MCP, $4\alpha_1 - 3\mu = 1.958$; for the Lasso, $4\alpha_1 - 3\mu = 2.108$.

We consider samples with sizes $T = 500, 1000, 1500, \dots, 20000$, and for each sample size, we simulate 100 times the process (\mathbf{Y}_t) . For each of these 100 simulated paths, we consider the sample variance covariance matrix $\hat{S} = \frac{1}{T} \sum_{t=1}^T \mathbf{X}_t \mathbf{X}_t^\top$. Thus, for each sample size, we obtain 100 LASSO/SCAD/MCP sparsity-based estimators $\hat{\Theta}^{\text{g}}, \hat{\Theta}^{\text{ls}}, \hat{\Theta}^{\text{vn}}, \hat{\Theta}^{\text{dt}}$. For the Stein (resp. von Neumann) case, we set $\eta = 1$ (resp. $\eta = 1$). The $\|\cdot\|_2$ (resp. $\|\cdot\|_1$) consistency patterns and the theoretical upper bounds for $\lambda_T = 4.5\sqrt{\log(d^2)/T}$ and $\lambda_T = 6.5\sqrt{\log(d^2)/T}$ are reported in Panels 1a-1f (resp. Panels 2a-2f). For all cases, the consistency is more favorable when using the rate $2\sqrt{\log(d^2)/T}$, which is in line with the theoretical upper bounds that depend on λ_T : using a tighter rate provides tighter bounds as depicted in our figures. In small samples, the least squares based estimator is less precise but performs well for large samples. Note that for the ℓ_1 -consistency and $\lambda_T = 4.5\sqrt{\log(d^2)/T}$, only the theoretical upper bound for the LASSO is reported: for all other cases, a large sample size $T \gg 20000$ would be required to reach from above $\|\text{vec}(\Theta_0)\|_1$. Panels 3a-3c report the consistency patterns for all loss based estimators for $\lambda_T = 4.5\sqrt{\log(d^2)/T}$. The von Neumann based estimators are close to each other since the $b_{\text{scad}}, b_{\text{mcp}}$ parameters are large so that the resulting penalization rate behaves as a LASSO one. A significant issue is how "informative" these error bounds are. Their rates depend on the regularization parameter λ_T , on the curvature of the loss function through the RSC parameters α_1 and the non-convexity of the penalty, where the trade-off expressed through the constraint $4\alpha_1 > 3\mu$ is key in our theoretical analysis. For the von Neumann loss, $4\alpha_1 > 3\mu$ is satisfied for significantly large values of b_{scad} and b_{mcp} so that μ is small enough compared to α_1 . Hence the denominator for the von Neumann case implies that the upper bounds for both the $\|\cdot\|_1$ and $\|\cdot\|_2$ errors are the least informative. For the ℓ_2 -consistency, it would require a dramatically large sample size T for the LASSO theoretical upper bound to cut the $\|\Theta_0\|_F$ line from above. As for the Stein's loss, it would require a sample size $T > 200000$ for the LASSO theoretical upper bound to cut the $\|\Theta_0\|_F$ line from above and

is thus not reported in Panels 1a-1f. This feature changes when considering the least squares and D-trace loss functions, where the RSC parameter α_1 is large enough so that the denominator becomes larger and the theoretical upper bounds become informative.

4.2 An illustration of the support recovery property

In this subsection, the number of zero and non-zero coefficients are expressed with respect to $\text{vech}(\Theta_0)$. We set $N = 30$ and $p = 2$ so that $d = 90$ and the total number of parameters in Θ is 8100 (idest 4095 distinct elements). The total number of zero components in $\text{vech}(\Theta)$ is 3676 so that the number of non-zero coefficients in $\text{vech}(\Theta)$ is 419. Here, $\|\text{vec}(\Theta_0)\|_1 = 141.35$, $\|\Theta_0\|_F = 10.28$, $\|\Theta_0\|_{\max} = 1.001$, $\|\Theta_0\|_s = 2.30$. Note that the total number of parameters in B_0 (resp. B_1, B_2 ; resp. $\text{vech}(\Sigma_u)$) is 788 (resp. 1664; resp. 394). Using the proposed sparse estimation to recover \mathcal{A} , for a sample size T , we simulated (20) a hundred times and assess the ability to correctly identify the support. We report in Table 1 the variable selection performance through the number of zero coefficients correctly estimated, denoted as C , the number of zero coefficients incorrectly estimated (i.e. an estimated zero coefficient whereas the true parameter is non-zero), denoted as IC1, the number of nonzero coefficients incorrectly estimated (i.e. an estimated non-zero coefficient whereas the true parameter is zero), denoted IC2, averaged for these hundred batches. The mean squared error is reported as an estimation accuracy measure. The Stein's (resp. von Neumann) loss is calibrated with $\eta = 1$ (resp. $\eta = 0$). The $b_{\text{scad}}, b_{\text{mcp}}$ are set so that $4\alpha_1 - 3\mu > 0$ is still satisfied: $b_{\text{scad}} = 50, b_{\text{mcp}} = 45$ (resp. 30, 20, resp. 150, 155, resp. 50, 45) for the Stein's (resp. Least squares, resp. von Neumann, resp. D-Trace) loss. The regularization parameter λ_T is set as $c\sqrt{\log(d^2)/T}$, with $c = 4$ a value calibrated by a cross-validation (CV) procedure and selected as optimal for the Stein's loss and D-Trace loss for $T = 30000$. For the sake of clarification, we set the same value for all losses. We used the data-dependent hv-CV procedure devised by Racine (2000), which consists in leaving a gap between the test sample and the training sample, on both sides of the test sample. Our simulation results show the challenge to perfectly recover the sparse model for all sparse estimators. First, Θ_0 contains a large number

of small non-zero coefficients, where $\min_{i \in \mathcal{A}} |\text{vech}(\Theta_0)| = 0.0022$ and the number of coefficients in absolute value smaller than 0.05 is 223. This mainly justifies the IC1 figures. For all cases, the LASSO provides higher MSE, which agrees with the property that the LASSO penalizes all coefficients with the same rate, thus generating a larger bias. We note that for larger sample sizes, the recovery becomes more precise since both IC1 and IC2 diminish.

4.3 Application to real data

To assess the relevance of the new estimation technique, we consider a VAR model based on a vector of squared stock returns to obtain forecasts using the three loss functions with the three penalties. We use the stocks listed in the Dow Jones Industrial Average, excluding Dow Inc. ($N = 29$) for the period starting from February 18, 2010 to June 19, 2018, providing 2100 observations. Dow Inc. is left aside since it went public on April 1st, 2019. The last 100 observations are used to carry out the forecasting analysis, where we apply a rolling window with sample size $T = 2000$. We selected the regularization parameters $\lambda_T = \sqrt{\log((p+1)^2 N^2)/T}$ for the Gaussian and the D-Trace loss functions and $\lambda_T = 0.2 \sqrt{\log((p+1)^2 N^2)/T}$ for the LS and von Neumann loss functions: to do so, we applied a cross validation over the 2000 first observations and then applied the same regularization rate over the next 100 forecasting periods. As a benchmark, we estimated the VAR model (1) with $p = 1, \dots, 5$ by the ordinary least squares (OLS) method. Table 2 presents the number of parameters in the model and AIC, indicating that the VAR(1) model has the minimum AIC.

We obtained the one-step-ahead forecasts for the VAR(1) model in order to provide the mean squared forecast error (MSFE) and the sample mean of the number of non-zero parameter estimates in B_0 and B_1 for the 100 forecasting period. 3 shows the forecasting results. The minimum number of non-zero parameters is represents less than 15 percent of the standard non-sparse OLS case. On the other hand, the sparse modelling of SVAR models have smaller MSFEs than the traditional VAR has. Among the proposed penalized losses, the von Neumann with SCAD has the minimum MSFE. Furthermore, to examine the

significance of the forecast loss, Table 3 presents p -value for the Model Confidence Set (MCS) of Hansen, Lunde and Nason (2011). The MCS procedure selects a set of models that contains the best model with 95%. Except for OLS and the Gaussian loss function with LASSO, all approaches are included in the same confidence set. Except for these LASSO cases, the new approach provides significantly superior results than the OLS does. The results for the real data analysis support the relevance of the method over the standard OLS method.

References

- Abadir, K.M. and J.R. Magnus (2005) *Matrix algebra*. Cambridge University Press.
- Ahelegbey, D. F., M. Billio and R. Casarin (2016) Bayesian graphical models for structural vector autoregressive processes. *Journal of Applied Econometrics*, **31**, 357–386.
- Alquier, P., K. Bertin, P. Doukhan and R. Garnier (2020) High dimensional VAR with low rank transition. *Statistics and Computing*, **30**, 1139–1153.
- Basu, S. and G. Michailidis (2015) Regularized estimation in sparse high-dimensional time series models. *The Annals of Statistics*, **43**, No. 4, 1535–1567.
- Bickel, P. J. and E. Levina (2008a) Covariance regularization by thresholding. *The Annals of Statistics*, **36**, 2577–2604.
- Bickel, P. J. and E. Levina (2008b) Regularized estimation of large covariance matrices. *The Annals of Statistics*, **36**, 199–227.
- Bien, J. and R. Tibshirani (2011) Sparse Estimation of a Covariance Matrix. *Biometrika*, **98**, 807–820.
- Blanchard, O. J. and D. Quah (1989) The dynamic effects of aggregate demand and supply disturbances. *American Economic Review*, **79**, 655–673.
- Candès, E.J. and Y. Plan (2009) Near-ideal model selection by ℓ_1 minimization. *The Annals of Statistics*, **37**, No. 5A, 2145–2177.
- Dahlhaus, R. and M. Eichler (2000) Causality and graphical models for time series. In P. Green, N. Hjort, and S. Richardson, eds., *Highly Structured Stochastic Systems*, Oxford University Press.
- Dedecker, J. and X. Fan (2015) Deviation inequalities for separately Lipschitz functionals of iterated random functions. *Stochastic Processes and their Applications*, **125**, No. 1, 60–90.
- Fan, J., Y. Feng and Y. Wu (2009) Network exploration via the adaptive lasso and scad penalties. *The Annals of Applied Statistics*, **3**, No. 2, 521–541.
- Fan, J. and R. Li (2001) Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, **96**, 1348–1360.
- Fan, J., Y. Liao and M. Mincheva (2011) High-dimensional covariance matrix estimation in approximate factor models. *The Annals of Statistics*, **39**, No. 6, 3320–3356
- Hansen, P. R., Lunde, A., and J. M. Nason (2011) The model confidence set. *Econometrica*, **79**, 453–497.
- Huang, J., N. Liu, M. Pourahmadi, and L. Liu (2006) Covariance matrix selection and estimation via penalized normal likelihood. *Biometrika*, **93**, 85–98.
- Inoue, A. and L. Killian (2013) Inference on impulse response functions in structural VAR models. *Journal of Econometrics*, **177**, 1–13.
- Johnston, J. (1972) *Econometric Methods*, 2nd ed., McGraw-Hill: New York.
- Lam, C. and J. Fan (2009) Sparsistency and rates of convergence in large covariance matrix estimation. *The Annals of Statistics*, **37**, 4254–4278.
- Loh, P.L. and M.J. Wainwright (2015) regularized M-estimators with non-convexity: statistical and algorithmic theory for local optima. *Journal of Machine Learning Research*, **16**, 559–616.
- Loh, P.L. and M.J. Wainwright (2017) Support recovery without incoherence: a case for non-convex regularization. *The Annals of Statistics*, **45**, No. 6, 2455–2482.
- Lütkepohl, H. (2006) *New Introduction to Multiple Time Series Analysis*, New York: Springer-Verlag.
- Lütkepohl, H. (2017) Estimation of structural vector autoregressive models. *Communications for Statistical Applications and Methods*, **24**, 421–441.
- Merlevède F., M. Peligrad and E. Rio (2011) A Bernstein type inequality and moderate deviations for weakly dependent sequences. *Probability Theory and Related Fields*, **151**, 435–474.
- Negahban, S.N, P. Ravikumar, M.J. Wainwright and B. Yu (2012) A unified framework for high-dimensional analysis of M-estimators with decomposable regularizers. *Statistical Science*, **27**, No. 4, 538–557.
- Oxley, L., M. Reale and G. Tunnicliffe-Wilson (2009) Constructing Structural VAR Models with Conditional Independence Graphs. *Mathematics and Computers in Simulation*, **79**, 2910–2916.

- Poignard, B. and J.D. Fermanian (2021) Finite sample properties of Sparse M-estimators with Pseudo-Observations. To appear in *Annals of the Institute of Statistical Mathematics*.
- Poignard, B. and Y. Terada (2020) Statistical Analysis of Sparse Approximate Factor Models. *Electronic Journal of Statistics*, **14**, 3315–3365.
- Racine, J. (2000). Consistent cross-validators model-selection for dependent data: hv-block cross-validation. *Journal of Econometrics*, **99**, 39–61.
- Ravikumar, P., M.J. Wainwright, G. Raskutti and B. Yu (2011) High-dimensional covariance estimation by minimizing ℓ_1 -penalized log-determinant divergence. *Electronic Journal of Statistics*, **5**, 935–980.
- Rothman, A. J., P. J. Bickel, E. Levina, and J. Zhu (2008) Sparse permutation invariant covariance estimation. *Electronic Journal of Statistics*, **2**, 494–515.
- Rothman, A. J., E. Levina, and J. Zhu (2009) Generalized thresholding of large covariance matrices. *Journal of the American Statistical Association*, **104**, 177–186.
- Rothman, A. J., E. Levina, and J. Zhu (2010) A new approach to Cholesky-based covariance regularization in high dimensions. *Biometrika*, **97**, 539–550.
- Sims, C. A. (1980) Macroeconomics and Reality. *Econometrica*, **48**, 1–48.
- Tunncliffe-Wilson, G. and M. Reale (2008) The sampling properties of conditional independence graphs for I(1) structural VAR models. *Journal of Time Series Analysis*, **29**, 802–810.
- Waggoner, D. F. and T. Zha (2003) Likelihood preserving normalization in multiple equation models. *Journal of Econometrics*, **114**, 329–347.
- Wainwright, M.J. (2009) Sharpe thresholds for high-dimensional and noisy sparsity recovery using ℓ_1 -constrained quadratic programming (Lasso). *IEEE Transactions on Information Theory*, **55**, No. 5, 2183–2202.
- White, H. (2001) *Asymptotic Theory for Econometricians*, 2nd edition, Emerald, UK
- Wong, K.C., Z. Li and A. Tewari (2020) Lasso guarantees for β -mixing heavy-tailed time series. *The Annals of Statistics*, **48**, No. 2, 1124–1142.
- Wu, Y. and T. Li (2020) Differential network inference via the fused D-trace loss with cross variables. *Electronic Journal of Statistics*, **14**, 1269–1301.
- Wu, W. B. and M. Pourahmadi (2003) Nonparametric estimation of large covariance matrices of longitudinal data. *Biometrika*, **90** 831–844.
- Yuan, M. and Y. Lin (2007) Model selection and estimation in the Gaussian graphical model. *Biometrika*, **94**, 19–35.
- Zhang, C.-H. (2010) Nearly unbiased variable selection under minimax concave penalty. *The Annals of Statistics*, **38**, 894–942.
- Zhang, T. and H. Zou (2014) Sparse precision matrix estimation via lasso penalized D-trace loss. *Biometrika*, **101**, 103–120.
- Zhao, P. and B. Yu (2006) On model selection consistency of Lasso. *Journal of Machine Learning Research*, **7**, 2541–2567.

5 Technical appendix

5.1 Figures and Tables

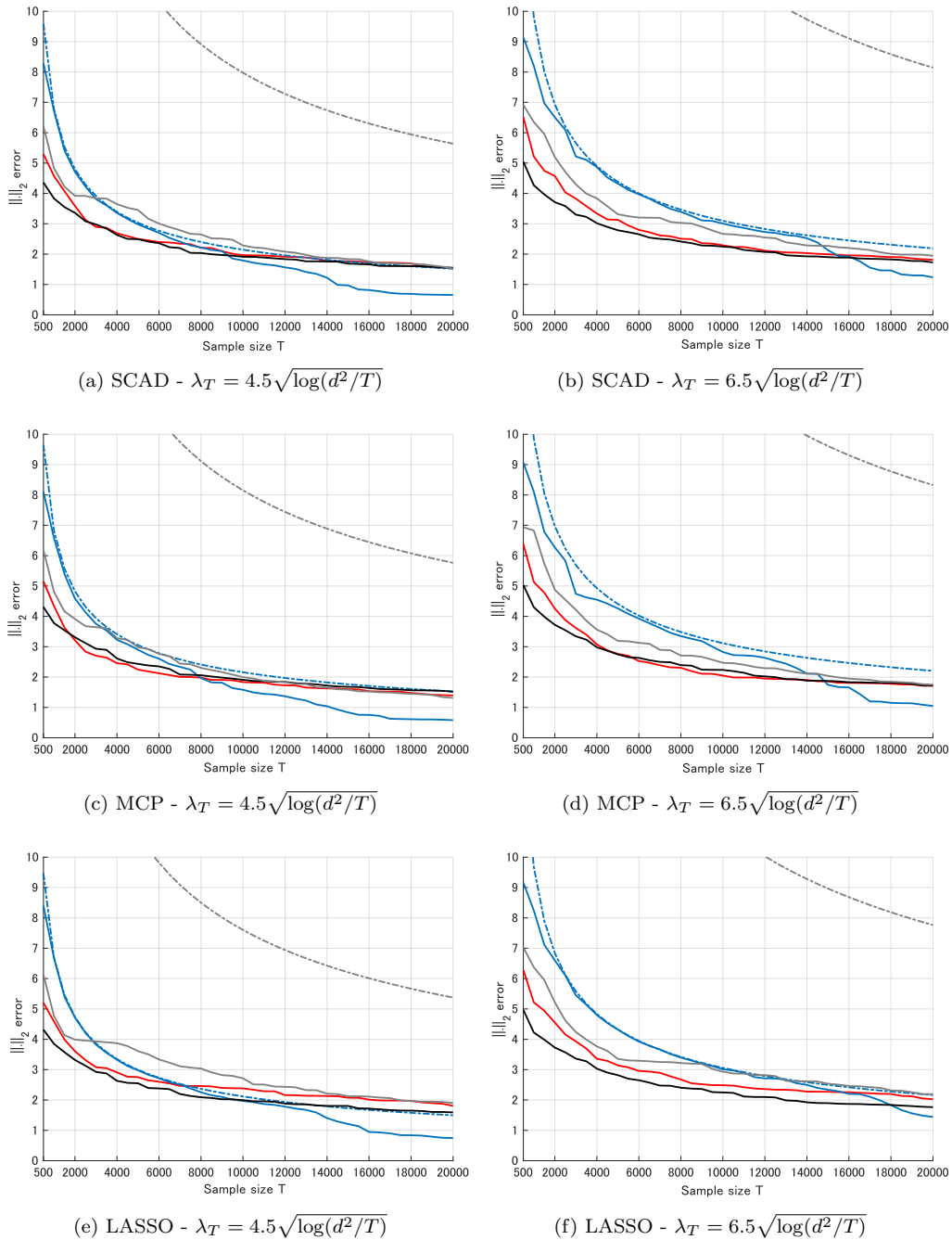


Figure 1: $\|\cdot\|_2$ consistency for the setting of Subsection 4.1 for $\lambda_T = 4.5\sqrt{\log(d^2)/T}$ and $\lambda_T = 6.5\sqrt{\log(d^2)/T}$. The Stein, least squares, von Neumann and D-trace cases are represented in solid red, blue, black and gray respectively. Each point represents an average of 200 trials for each sample size. The corresponding theoretical upper bounds are in dashed-dotted lines.

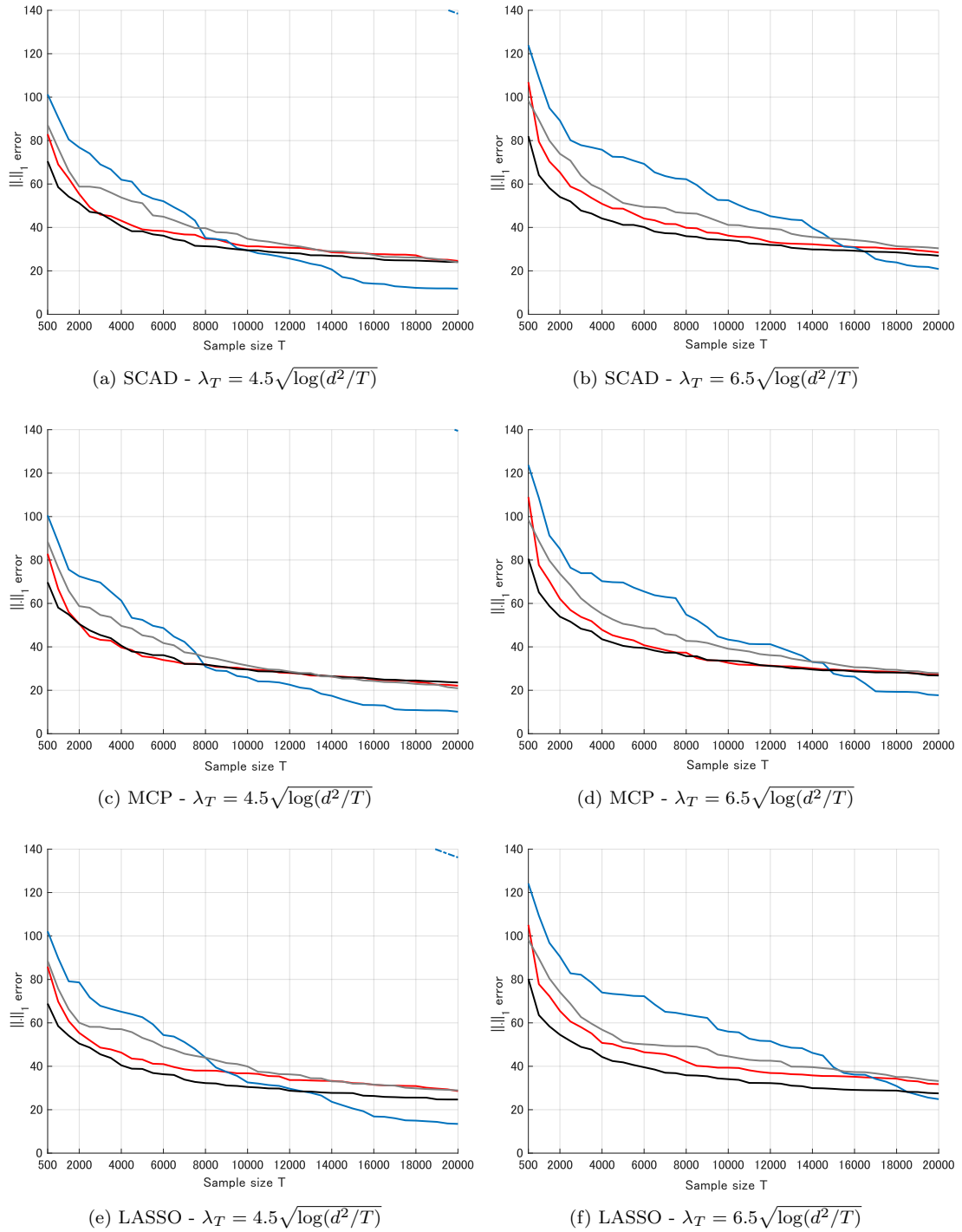
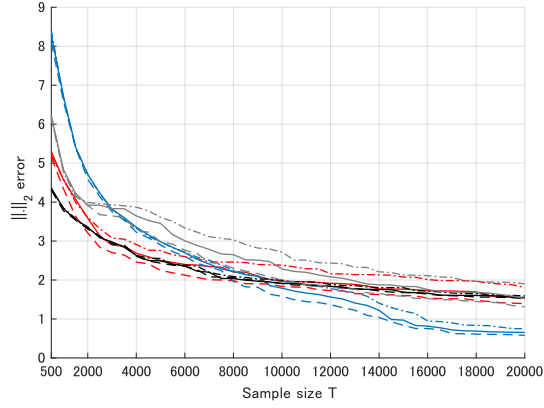
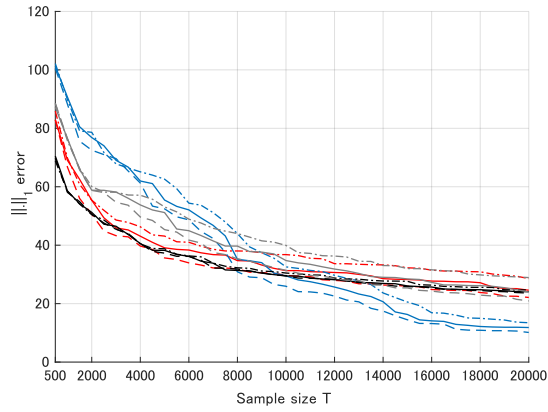


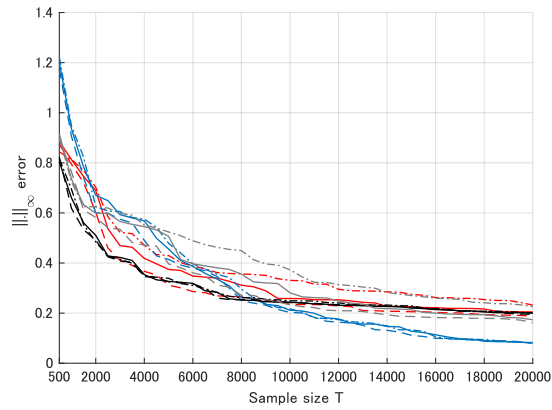
Figure 2: $\|\cdot\|_1$ consistency for the setting of Subsection 4.1 for $\lambda_T = 4.5\sqrt{\log(d^2)/T}$ and $\lambda_T = 6.5\sqrt{\log(d^2)/T}$. The Stein, least squares, von Neumann and D-trace cases are represented in solid red, blue, black and gray respectively. Each point represents an average of 200 trials for each sample size. The corresponding theoretical upper bounds are in dashed-dotted lines.



(a) $\|\cdot\|_2$ -consistency



(b) $\|\cdot\|_1$ -consistency



(c) $\|\cdot\|_\infty$ -consistency

Figure 3: $\|\cdot\|_2, \|\cdot\|_1, \|\cdot\|_\infty$ consistencies for the setting of Subsection 4.1 for $\lambda_T = 4.5\sqrt{\log(d^2)/T}$. The Stein, least squares, von Neumann and D-trace cases are represented in red, blue, black and gray color respectively. For each loss, the SCAD (resp. MCP, resp. LASSO) is in solid (resp. dashed, resp. dashed-dotted) line. Each point represents an average of 200 trials for each sample size.

Table 1: Model selection and precision accuracy for $N = 30$ and $p = 2$ based on 100 replications. For each penalized Stein, Least squares, Von Neumann and D-trace, the penalty cases are reported according to the order: SCAD, MCP and LASSO

	Truth	Stein	Least-squares	von Neumann	D-trace
$T = 5000$					
C	3676	3469.1 – 3458.8 – 3464.8	3535.4 – 3523.2 – 3535.1	3432.9 – 3422.5 – 3426.6	3462.6 – 3463.6 – 3462.3
IC1	0	222.5 – 221.3 – 222.5	217.8 – 218.6 – 218.5	222.6 – 220.9 – 221.8	223.7 – 222.2 – 224.9
IC2	0	206.9 – 217.1 – 211.2	140.6 – 152.8 – 140.9	243.1 – 253.5 – 249.4	213.4 – 212.7 – 213.7
MSE		3.291 – 3.342 – 3.549	7.085 – 7.522 – 7.058	1.462 – 1.459 – 1.418	3.267 – 3.347 – 3.598
$T = 10000$					
C	3676	3603.4 – 3610.2 – 3609.9	3602.1 – 3601.9 – 3607.7	3583.5 – 3574.2 – 3579.8	3612.6 – 3601.8 – 3602.1
IC1	0	216.1 – 214.7 – 217.2	210.3 – 206.5 – 219.6	215.9 – 211.7 – 220.8	217.7 – 208.1 – 221.9
IC2	0	72.5 – 65.8 – 66.0	73.9 – 74.1 – 68.3	92.5 – 101.7 – 96.2	63.4 – 74.2 – 73.9
MSE		2.332 – 2.341 – 2.787	5.663 – 6.014 – 6.405	0.970 – 0.967 – 1.211	1.951 – 1.955 – 2.329
$T = 20000$					
C	3676	3643.5 – 3643.6 – 3650.3	3639.1 – 3637.4 – 3644.3	3649.9 – 3650.9 – 3648.7	3653.5 – 3654.2 – 3657.3
IC1	0	192.4 – 183.8 – 203.7	193.6 – 180.2 – 196.1	189.3 – 179.8 – 196.2	184.2 – 171.8 – 193.6
IC2	0	32.4 – 32.3 – 35.7	36.9 – 38.6 – 41.6	26.1 – 25.0 – 27.3	22.5 – 21.8 – 28.6
MSE		1.446 – 1.491 – 1.964	2.554 – 2.514 – 3.242	0.643 – 0.637 – 0.758	0.961 – 1.012 – 1.391
$T = 30000$					
C	3676	3665.1 – 3670.2 – 3674.9	3668.6 – 3673.4 – 3674.4	3675.7 – 3675.5 – 3675.6	3665.1 – 3665.9 – 3675.8
IC1	0	172.8 – 161.8 – 188.2	187.7 – 170.2 – 188.9	182.3 – 171.1 – 191.3	178.1 – 163.4 – 186.7
IC2	0	16.9 – 17.7 – 19.1	7.4 – 2.5 – 1.6	0.3 – 0.5 – 0.8	10.9 – 7.3 – 0.2
MSE		0.185 – 0.226 – 0.396	0.254 – 0.156 – 0.293	0.113 – 0.110 – 0.163	0.102 – 0.117 – 0.202

Table 2: OLS Results

VAR(p)	$p = 1$	$p = 2$	$p = 3$	$p = 4$	$p = 5$
# of parameters	1276	2117	2958	3799	4640
AIC	16673.7	16841.8	17009.9	17178.1	17346.2

Table 3: Sparse VAR Forecasts

Loss	Penalty	# of Non-Zero	MSFE	MCS
Gaussian	SCAD	180.0	448.01	0.080
Gaussian	MCP	181.8	443.88	0.080
Gaussian	LASSO	180.5	451.76	0.035
LS	SCAD	479.0	442.37	0.080
LS	MCP	477.3	441.44	0.080
LS	LASSO	488.1	441.77	0.080
D-Trace	SCAD	359.7	435.42	0.454
D-Trace	MCP	356.7	436.31	0.080
D-Trace	LASSO	359.7	435.27	0.454
von Neumann	SCAD	296.6	431.94	1.000
von Neumann	MCP	289.9	432.63	0.454
von Neumann	LASSO	291.9	432.02	0.911
OLS	—	1276	457.65	0.035

Note: The results are based on the one-step-ahead forecasts for the last 100 observations using the rolling window with the sample size $T = 2000$. ‘# of Non-Zero’ indicates the sample mean of the number of non-zero estimates for B_0 and B_1 in the (sparse) VAR(1) model. ‘MSFE’ denotes the mean squared forecasts error. The column labeled MCS presents the p -values associated with the Model Confidence Set of Hansen et al. (2011).

5.2 Intermediary results

We provide the primal dual witness method as in Loh and Wainwright (2017), an approach that relies on the following steps. Here, the parameter of interest is $\theta \in \mathbb{R}^q$ and q denotes the dimension. The problem of interest is a regularized M-estimation one, where a generic loss $\mathbb{G}_T(\cdot)$ is penalized by $\mathbf{p}(\lambda_T, \cdot)$:

$$\hat{\theta} = \arg \min_{\theta \in \Omega} \left\{ \mathbb{G}_T(\theta) + \mathbf{p}(\lambda_T, \theta) \right\}, \quad \Omega = \left\{ \theta \in \Theta \subseteq \mathbb{R}^q, \|\theta\|_1 \leq R \right\}. \quad (21)$$

The loss function $\mathbb{G}_T(\cdot)$ satisfies the RSC condition and $\mathbf{p}(\lambda_T, \cdot)$ assumption 2.

Step 1. We define the estimator

$$\hat{\theta}_{\mathcal{A}} = \arg \min_{\theta \in \mathbb{R}^{|\mathcal{A}|}: g(\theta) \leq R, \theta \in \Omega} \left\{ \mathbb{G}_T(\theta) + \mathbf{p}(\lambda_T, \theta) \right\}. \quad (22)$$

We solve problem (22), under the constraint $\hat{\mathcal{A}} \subseteq \mathcal{A}$ and prove $\|\hat{\theta}_{\mathcal{A}}\|_1 < R$.

Step 2. Defining $\hat{\mathbf{z}}_{\mathcal{A}} \in \partial \|\hat{\theta}_{\mathcal{A}}\|$, we choose $\hat{\mathbf{z}}_{\mathcal{A}^c}$ satisfying the orthogonality condition

$$\nabla_{\theta} \mathbb{G}_T(\hat{\theta}) - \nabla_{\theta} \mathbf{q}(\lambda_T, \hat{\theta}) + \lambda_T \hat{\mathbf{z}} = 0, \quad (23)$$

with $\hat{\mathbf{z}} = (\hat{\mathbf{z}}_{\mathcal{A}}, \hat{\mathbf{z}}_{\mathcal{A}^c})$, $\hat{\theta} = (\hat{\theta}_{\mathcal{A}}, 0_{\mathcal{A}^c})$, $\mathbf{q}(\lambda_T, \rho) = \lambda_T \rho - \mathbf{p}(\lambda_T, \rho)$. We prove strict dual feasibility $\|\hat{\mathbf{z}}_{\mathcal{A}^c}\|_{\infty} < 1$.

Step 3. We prove that $\hat{\theta}$ is a local optimum of (21) and that any stationary point of (21) satisfies $\text{supp}(\hat{\theta}) \subseteq \mathcal{A}$.

The PDW procedure does not allow for *practically* solving the regularization problem (21) as step 1 requires to know the true subset model \mathcal{A} . However, this approach is useful as a proof method to characterize the optimal solution $\hat{\theta}$. In **Step 1**, the criterion (22) is strictly convex under the RSC condition. This implies that for $\|\hat{\theta}_{\mathcal{A}}\|_1 < R$, the subgradient condition (23) must hold at $\hat{\theta}_{\mathcal{A}}$ for the restricted problem (22). Loh and Wainwright (2017) proves that, although problem (22) may be non-convex, the RSC condition and regularity conditions on the penalty function allow them to prove that the optimum obtained in **Step 3** is a local optimum: see in particular their Lemma 10.

Using optimization reasoning, Loh and Wainwright (2017) provide conditions on λ_T, R to ensure the success of the PDW technique, which depends on **Step 3**, under the assumption that $\mathbb{G}_T(\cdot)$ satisfies the RSC condition with parameters $(\alpha_k, \tau_k)_{k=1,2}$ and $4\alpha_1 > 3\mu$. Indeed, these conditions guarantee that the support of $\hat{\theta}$ satisfying (23) in **Step 2** is the unique stationary point of the criterion (21): to be precise, the first condition concerns the suitable scaling of λ_T and R ; the second condition ensures strict dual feasibility - that is $\|\hat{\mathbf{z}}_{\mathcal{A}^c}\|_\infty < 1$ in **Step 2**. This is the object of the following Theorem.

Theorem 5.1. Loh and Wainwright (Theorem 1, 2017)

Suppose $\mathbb{G}_T(\cdot)$ satisfies the RSC condition with $(\alpha_k, \tau_k)_{k=1,2}$ parameters and $\mathbf{p}(\lambda_T, \cdot)$ is a μ -amenable penalty, with $0 \leq \mu < \alpha_1$. Suppose

(i) The parameters (λ_T, R) satisfy

$$4 \max \left\{ \|\nabla_\theta \mathbb{G}_T(\theta_0)\|_\infty, \alpha_2 \sqrt{\frac{\log(k_0)}{T}} \right\} \leq \lambda_T \leq \sqrt{\frac{(4\alpha_1 - 3\mu)\alpha_2}{384k_0}}, \quad (24)$$

$$\max \left\{ 2\|\theta_0\|_1, \frac{48k_0\lambda_T}{4\alpha_1 - 3\mu} \right\} \leq R \leq \min \left\{ \frac{\alpha_2}{8\lambda_T}, \frac{\alpha_2}{\tau_2} \sqrt{\frac{T}{\log(q)}} \right\}. \quad (25)$$

(ii) For some $\delta \in \left[\frac{4R\tau_1 \log(q)}{T\lambda_T}, 1 \right]$, the vector $\hat{\mathbf{z}}$ from the PDW construction satisfies the strict dual feasibility condition

$$\|\hat{\mathbf{z}}_{\mathcal{A}^c}\|_\infty \leq 1 - \delta. \quad (26)$$

Then for any k_0 -sparse vector θ_0 , the program (21) with a sample size $T \geq \frac{2\tau_1}{2\alpha_1 - \mu} k_0 \log q$ has a unique stationary point given by the primal output $\hat{\theta}$ of the PDW construction.

Suitable calibrations of λ_T and R , and thus proper scaling behaviours (T, q, k_0) , are necessary. Using exponential bounds, it is possible to evaluate the probability of satisfying (24) and (25) and thus the probability of the PDW success. In all their applications of interest - linear model, generalized linear model, Gaussian graphical Lasso - Loh and Wainwright (2017) obtain the upper bound $\|\nabla_\theta \mathbb{G}_T(\theta_0)\|_\infty \leq C\sqrt{\log(q)/T}$ with high probability. This motivates the choice λ_T proportional to $\sqrt{\log(q)/T}$ to satisfy

(24). Finally, it is worth noting that the trade-off between the curvature of the loss function through α_1 and the non-convexity degree of the penalty function through μ appears. As our simulations emphasize this trade-off for the Stein or von Neumann loss in particular, significantly large values for $b_{\text{scad}}, b_{\text{mcp}}$ are necessary to ensure $4\alpha_1 > 3\mu$.

In their Theorem 2, Loh and Wainwright (2017) provide an additional error bound under the conditions of Theorem 5.1. It also provides the guarantees that the unique optimum - local or global - (21) is the oracle estimator. The latter is defined as the non-penalized estimator obtained from minimizing the criterion $\mathbb{G}_T(\theta)$ over the true support \mathcal{A} . This is the object of the following Theorem.

Theorem 5.2. Loh and Wainwright (Theorem 2, 2017)

Under the conditions of Theorem 5.1, suppose strict dual feasibility (26) holds, suppose $\mathbf{p}(\lambda_T, \cdot)$ is μ -amenable with $\mu \in [0, \alpha_1)$. Then the unique stationary solution of (21) satisfies

(i)

$$\|\widehat{\theta} - \theta_0\|_\infty \leq \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{G}_T(\theta_0)_{\mathcal{A}}\|_\infty + \lambda_T \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty,$$

$$\text{with } \widehat{\mathbf{K}} = \int_0^1 \nabla_{\theta\theta^\top}^2 \mathbb{G}_T(\theta_0 + u(\widehat{\theta} - \theta_0)) du.$$

(ii) *If $\mathbf{p}(\lambda_T, \cdot)$ is (μ, ζ) -amenable and if the lower bound*

$$\min_{i \in \mathcal{A}} |\theta_{0,i}| \geq \lambda_T (\zeta + \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty) + \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{G}_T(\theta_0)_{\mathcal{A}}\|_\infty,$$

holds, then $\widehat{\theta}$ agrees with the oracle estimator $\widehat{\theta}^\circ$ and we have the bound

$$\|\widehat{\theta} - \theta_0\|_\infty \leq \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{G}_T(\theta_0)_{\mathcal{A}}\|_\infty.$$

These inequalities are expressed in a deterministic manner. As in Theorem 5.1, exponential bounds allow for upper bounding $\|\nabla_\theta \mathbb{G}_T(\theta_0)_{\mathcal{A}}\|_\infty$, which will provide explicit convergence rates over the $\|\cdot\|_\infty$ -error. The application of Theorem 5.2 requires that strict dual feasibility holds under the RSC condition. In their Proposition 1, Loh and Wainwright (2017) provide sufficient conditions to satisfy strict dual feasibility

for (μ, ζ) -amenable penalties, which thus allows for using Theorem 5.2. These conditions are given in the following Proposition.

Proposition 5.3. Loh and Wainwright (Proposition 1, 2017)

Under the conditions of Theorem 5.1, suppose $\mathbf{p}(\lambda_T, \cdot)$ is (μ, ζ) -amenable. Suppose

$$\theta_{0,\min} \geq \lambda_T(\zeta + \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty) + \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{G}_T(\theta_0)_{\mathcal{A}}\|_\infty,$$

with $\theta_{0,\min} = \min_{i \in \mathcal{A}} |\theta_i|$ and $\widehat{\mathbf{K}} = \int_0^1 \nabla_{\theta\theta^\top}^2 \mathbb{G}_T(\theta_0 + u(\widehat{\theta} - \theta_0)) du$. Then strict dual feasibility holds provided

$$\|\nabla_\theta \mathbb{G}_T(\theta_0)\|_\infty \leq \frac{1-\delta}{2} \lambda_T, \quad \text{and} \quad (27)$$

$$\|\widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{G}_T(\theta_0)_{\mathcal{A}}\|_\infty \leq \frac{1-\delta}{2} \lambda_T. \quad (28)$$

6 Proofs

Proof of Lemma 3.2. Under assumptions 3, 4 and 5, by Lemma A.2 of Fan, Liao and Mincheva (2013), $\mathbf{X}_{i,t} \mathbf{X}_{j,t}$ satisfies the exponential tail condition with parameter $\gamma_2/3$. Thus, by Theorem 1 of Merlevède, Peligrad and Rio (2011), there exist constants C_1, C_2, C_3, C_4, C_5 depending only on c, γ and γ_1 such that, for any positive ϵ and any $i, j \leq d^2$,

$$\begin{aligned} & \mathbb{P}\left(\left|\frac{1}{T} \sum_{t=1}^T \mathbf{X}_{i,t} \mathbf{X}_{j,t} - \Sigma_{x,ij}\right| > \epsilon\right) \\ & \leq T \exp\left(-\frac{(T\epsilon)^\gamma}{C_1}\right) + \exp\left(-\frac{(T\epsilon)^2}{C_2(1+TC_3)}\right) + \exp\left(-\frac{(T\epsilon)^2}{C_4 T} \exp\left(\frac{(T\epsilon)^{\gamma(1-\gamma)}}{C_5(\log(T\epsilon))^\gamma}\right)\right), \end{aligned}$$

so that $\mathbb{P}(\|\widehat{\mathbf{S}} - \Sigma_x\|_{\max} > \epsilon) = \mathbb{P}(\max_{1 \leq i, j \leq d} |\widehat{\mathbf{S}}_{ij} - \Sigma_{x,ij}| > \epsilon) \leq d^2 \max_{1 \leq i, j \leq d} \mathbb{P}(|\widehat{\mathbf{S}}_{ij} - \Sigma_{x,ij}| > \epsilon)$. We then deduce

$$\begin{aligned} & \mathbb{P}\left(\left|\frac{1}{T} \sum_{t=1}^T \mathbf{X}_{i,t} \mathbf{X}_{j,t} - \Sigma_{x,ij}\right| > \epsilon\right) \\ & \leq d^2 \left\{ T \exp\left(-\frac{(T\epsilon)^\gamma}{C_1}\right) + \exp\left(-\frac{(T\epsilon)^2}{C_2(1+TC_3)}\right) + \exp\left(-\frac{(T\epsilon)^2}{C_4 T} \exp\left(\frac{(T\epsilon)^{\gamma(1-\gamma)}}{C_5(\log(T\epsilon))^\gamma}\right)\right) \right\}. \end{aligned}$$

Let us now fix $\epsilon = L\sqrt{\log(d^2)/T}$ with $L > 0$ a constant large enough. For the scaling behaviour $T \geq K \log(d^2)^{2/\gamma-1}$ for a sufficiently large and positive constant K , then

$$\begin{aligned} & d^2 T \exp\left(-\frac{(T\epsilon)^\gamma}{C_1}\right) + d^2 \exp\left(-\frac{(T\epsilon)^2}{C_4 T} \exp\left(\frac{(T\epsilon)^{\gamma(1-\gamma)}}{C_5 (\log(T\epsilon))^\gamma}\right)\right) \\ &= \exp\left(-\frac{T^\gamma L^\gamma (\sqrt{\log(d^2)/T})^\gamma}{C_1} + \log(d^2) + \log(T)\right) + \exp\left(-\frac{T^2 \log(d^2)/T}{C_4 T} \exp\left(\frac{T^{\gamma(1-\gamma)} (L\sqrt{\log(d^2)/T})^{\gamma(1-\gamma)}}{C_5 (\log(TL\sqrt{\log(d^2)/T}))^\gamma}\right)\right) \\ &= o(\exp(-\log(d^2))), \end{aligned}$$

and

$$d^2 \exp\left(-\frac{(T\epsilon)^2}{C_2(1+TC_3)}\right) = \exp\left(-\frac{T^2(L\sqrt{\log(d^2)/T})^2}{C_2(1+TC_3)} + \log(d^2)\right) = O(\exp(-\log(d^2))).$$

As a consequence, if we fix $\epsilon = L\sqrt{\log(d^2)/T}$, then $\|\hat{S} - \Sigma_x\|_{\max} \leq L\sqrt{\log(d^2)/T}$ with probability at least $1 - O(\exp(-\log(d^2))) - o(\exp(-\log(d^2)))$. \square

Proof of Corollary 3.3. We first establish the RSC property. To do so, we derive the first and second order derivatives of the Stein's loss criterion $\mathbb{L}_{T,\eta}(\Theta)$ defined in (5). Using the differential operator applied with respect to Θ , we obtain

$$d\mathbb{L}_{T,\eta}(\Theta) = (1-2\eta)\Theta^{-1}(d\Theta) + \text{tr}\left(\eta\hat{S}(d\Theta) - (1-\eta)\Theta^{-1}(d\Theta)\Theta^{-1}\hat{S}^{-1}\right).$$

Hence in vector and matrix forms, the derivatives become

$$\nabla_{\theta}\mathbb{L}_{T,\eta}(\Theta) = \text{vec}\left((1-2\eta)\Theta^{-1} + \eta\hat{S} - (1-\eta)\Theta^{-1}\hat{S}^{-1}\Theta^{-1}\right), \quad \nabla_{\Theta}\mathbb{L}_{T,\eta}(\Theta) = (1-2\eta)\Theta^{-1} + \eta\hat{S} - (1-\eta)\Theta^{-1}\hat{S}^{-1}\Theta^{-1}.$$

Note that alternatively, the $\text{vech}(\cdot)$ operator could be applied to treat the redundant terms. Taking the $\|\cdot\|_{\infty}$ norm, over Θ_0 , the gradient becomes

$$\|\nabla_{\theta}\mathbb{L}_{T}(\Theta_0)\|_{\infty} = \|\nabla_{\Theta}\mathbb{L}_{T}(\Theta_0)\|_{\max} \leq \|(1-2\eta)\Theta^{-1} + \eta\hat{S} - (1-\eta)\Theta^{-1}\hat{S}^{-1}\Theta^{-1}\|_{\max}.$$

We now focus on the Hessian matrix. The second order differential is given by

$$d^2\mathbb{L}_{T,\eta}(\Theta) = (2\eta-1)\text{tr}\left(\Theta^{-1}(d\Theta)\Theta^{-1}(d\Theta)\right) + (1-\eta)\text{tr}\left(\Theta^{-1}(d\Theta)\Theta^{-1}\hat{S}^{-1}\Theta^{-1}(d\Theta) + \Theta^{-1}\hat{S}^{-1}\Theta^{-1}(d\Theta)\Theta^{-1}(d\Theta)\right).$$

We aim at extracting the form $\text{tr}(L(\mathbf{d}\Lambda)^\top M(\mathbf{d}\Lambda))$ for L (resp. M) any square $m \times m$ matrix (resp. $p \times p$).

We have

$$\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,\eta}(\Theta) = (2\eta - 1) \left(\Theta^{-1} \otimes \Theta^{-1} \right) + (1 - \eta) \left\{ \Theta^{-1} \otimes \Theta^{-1} \widehat{S}^{-1} \Theta^{-1} + \Theta^{-1} \widehat{S}^{-1} \Theta^{-1} \otimes \Theta^{-1} \right\}.$$

Now for some $\Theta_1 \in \Omega$ and $u \in [0, 1]$, let us define $\Theta = \Theta_0 + u\Gamma$ with $\Gamma = \Theta_1 - \Theta_0$. Then $\Theta \in \Omega$ and

$$\begin{aligned} f_T(\Theta) &:= \text{vec}(\Gamma)^\top \nabla_{\text{vec}(\Theta)\text{vec}(\Theta)^\top}^2 \mathbb{L}_{T,\eta}(\Theta) \text{vec}(\Gamma) \\ &\geq \text{vec}(\Gamma)^\top \left[(2\eta - 1) \left(\Theta^{-1} \otimes \Theta^{-1} \right) + (1 - \eta) \left\{ \Theta^{-1} \otimes \Theta^{-1} \widehat{S}^{-1} \Theta^{-1} + \Theta^{-1} \widehat{S}^{-1} \Theta^{-1} \otimes \Theta^{-1} \right\} \right] \text{vec}(\Gamma) \\ &\geq \|\Gamma\|_F^2 \left[(2\eta - 1) \lambda_{\min}(\Theta^{-1} \otimes \Theta^{-1}) + (1 - \eta) \lambda_{\min} \left(\left\{ \Theta^{-1} \otimes \Theta^{-1} \widehat{S}^{-1} \Theta^{-1} + \Theta^{-1} \widehat{S}^{-1} \Theta^{-1} \otimes \Theta^{-1} \right\} \right) \right] \\ &\geq \|\Gamma\|_F^2 \left[(2\eta - 1) \lambda_{\min}(\Theta^{-1})^2 + 2(1 - \eta) \lambda_{\min}(\Theta^{-1})^3 \lambda_{\min}(\widehat{S}^{-1}) \right] \end{aligned}$$

because the spectrum of $A \otimes B$ is the cross product of the spectrums of A and B (see, e.g., Lütkepohl, 1996, Subsection 5.2.1), and $\lambda_{\min}(\Theta) = \inf_{\mathbf{x}} \mathbf{x}^\top \Theta \mathbf{x} / \|\mathbf{x}\|_2^2$. We now focus on $\lambda_{\min}(\Theta^{-1})^2$. We have $\lambda_{\max}(\Theta) \leq \lambda_{\max}(\Theta_0) + \lambda_{\max}(\Theta_1 - \Theta_0)$, which implies $\lambda_{\min}(\Theta^{-1})^2 \geq \{\lambda_{\max}(\Theta_0) + 1\}^{-2}$. Therefore

$$f_T(\Theta) \geq \|\Gamma\|_F^2 \left[(2\eta - 1) \{\lambda_{\max}(\Theta_0) + 1\}^{-2} + 2(1 - \eta) \{\lambda_{\max}(\Theta_0) + 1\}^{-3} \lambda_{\min}(\widehat{S}^{-1}) \right].$$

The true vector of parameters is $\theta_0 = \text{vec}(\Theta_0)$ and we have $\|\Gamma\|_F^2 = \|\Theta - \Theta_0\|_2^2$. As a consequence,

$$(\theta - \theta_0)^\top \nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\bar{\Theta})(\theta - \theta_0) \geq \|\theta - \theta_0\|_2^2 \left[(2\eta - 1) \{\lambda_{\max}(\Theta_0) + 1\}^{-2} + 2(1 - \eta) \{\lambda_{\max}(\Theta_0) + 1\}^{-3} \lambda_{\min}(\widehat{S}^{-1}) \right],$$

for any $\bar{\Theta}$ that lies between Θ and Θ_0 . Thus, at Θ_0 , the RSC condition is satisfied with coefficients $\alpha_1 = (2\eta - 1) \{\lambda_{\max}(\Theta_0) + 1\}^{-2} + 2(1 - \eta) \{\lambda_{\max}(\Theta_0) + 1\}^{-3} \lambda_{\min}(\widehat{S}^{-1})$ and $\alpha_2 = \alpha_1$, $\tau_1 = \tau_2 = 0$. Should we take $\eta = 1$, then $\alpha_1 = \{\lambda_{\min}(\Theta_0) + 1\}^{-2}$ since the Hessian simplifies as $\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,\eta}(\Theta) = \Theta^{-1} \otimes \Theta^{-1}$.

We now evaluate the probability of satisfying (4) when $\eta = 1$, a situation most commonly used when dealing with inverse variance covariance estimation. Then:

$$\|\nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)\|_\infty = \|\nabla_{\Theta} \mathbb{L}_{T,1}(\Theta_0)\|_{\max} = \|\widehat{S} - \Sigma_x\|_{\max},$$

which provides (6). By Lemma 3.2, we conclude $\|\nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)\|_\infty \leq K \sqrt{\frac{\log(d)}{T}}$, with probability at least $1 - O(\exp(-\log(d^2))) - o(\exp(-\log(d^2)))$ for a sample size $T \geq L \log(d^2)^{2/\gamma-1}$. \square

Proof of Corollary 3.4. We first establish the RSC condition for the least squares loss function. Using the differential operator with respect to Θ , we have $d\mathbb{L}_T(\Theta) = 2\text{tr}(\Theta - \widehat{S}^{-1})(d\Theta)$. Hence

$$\nabla_{\theta}\mathbb{L}_T(\Theta) = 2\text{vec}(\Theta - \widehat{S}^{-1}), \quad \nabla_{\Theta}\mathbb{L}_T(\Theta) = 2(\Theta - \widehat{S}^{-1}).$$

As for the Hessian, by identification, we deduce

$$\nabla_{\theta\theta^\top}^2\mathbb{L}_T(\Theta) = 2(I_d \otimes I_d).$$

We thus deduce $(\theta - \theta_0)^\top \nabla_{\theta\theta^\top}^2\mathbb{L}_T(\bar{\Theta})(\theta - \theta_0) \geq 2\|\theta - \theta_0\|_2^2$, for any $\bar{\Theta}$ that lies between Θ and Θ_0 . Thus, at Θ_0 , the RSC condition is satisfied with coefficients $\alpha_1 = 2$ and $\alpha_2 = \alpha_1$, $\tau_1 = \tau_2 = 0$.

We now evaluate the probability of satisfying (4). We have

$$\|\nabla_{\theta}\mathbb{L}_T(\Theta_0)\|_{\infty} = \|\nabla_{\Theta}\mathbb{L}_T(\Theta_0)\|_{\max} = 2\|\Theta_0 - \widehat{S}^{-1}\|_{\max} \leq 2\|\Theta_0 - \widehat{S}^{-1}\|_s.$$

Now, with high probability, for a sufficiently large constant K , $\|\widehat{S} - \Sigma_x\|_{\max} \leq K\sqrt{\frac{\log(d^2)}{T}}$, and we aim at bounding $\|\widehat{S}^{-1} - \Sigma_x^{-1}\|_{\max}$. To do so, we use series expansion for the inverse of matrices. By the Woodbury matrix identities, we have for two symmetric positive definite matrices A and B :

$$(A - B)^{-1} = \sum_{k=0}^{\infty} (A^{-1}B)^k A^{-1} = A^{-1} + \sum_{k=1}^{\infty} (A^{-1}B)^k A^{-1}.$$

Thus

$$(A + B)^{-1} = A^{-1} + \sum_{k=1}^{\infty} (-1)^k (A^{-1}B)^k A^{-1}.$$

Now, taking $B = \widehat{S} - \Sigma_x$ and $A = \Sigma_x$, then we obtain

$$\widehat{S}^{-1} - \Sigma_x^{-1} = \sum_{k=1}^{\infty} (-1)^k (\Sigma_x^{-1}(\widehat{S} - \Sigma_x))^k \Sigma_x^{-1}.$$

Hence, using this relationship, we can upper bound the norm of the difference $\widehat{S}^{-1} - \Sigma_x^{-1}$ by the norm of $\widehat{S} - \Sigma_x^{-1}$ for which we have an exponential bound. For any sub-multiplicative matrix norm $\|\cdot\|$, we deduce

$$\|\widehat{S}^{-1} - \Sigma_x^{-1}\| \leq \|\Sigma_x^{-1}\| \sum_{k=1}^{\infty} \left(\|\Sigma_x^{-1}\| \|\widehat{S} - \Sigma_x\| \right)^k = \|\Sigma_x^{-1}\|^2 \|\widehat{S} - \Sigma_x\| \sum_{k=0}^{\infty} \left(\|\Sigma_x^{-1}\| \|\widehat{S} - \Sigma_x\| \right)^k.$$

For a square symmetric matrix M , under the condition $\|M\|_s < 1$, then $\sum_{k=0}^{\infty} \|M\|_s^k = (1 - \|M\|_s)^{-1}$. Now, with high probability, for a sufficiently large sample size, we thus have

$$\|\widehat{S}^{-1} - \Sigma_x^{-1}\|_s \leq \frac{\|\Sigma_x^{-1}\|_s^2 \|\widehat{S} - \Sigma_x\|_s}{1 - \|\Sigma_x^{-1}\|_s \|\widehat{S} - \Sigma_x\|_s},$$

so that bounding $\|\widehat{S} - \Sigma_x\|_s$ implies bounding $\|\widehat{S}^{-1} - \Sigma_x^{-1}\|_s$. Now, we have

$$\|\Sigma_x^{-1} - \widehat{S}^{-1}\|_{\max} \leq \|\Sigma_x^{-1} - \widehat{S}^{-1}\|_s \leq \frac{\|\widehat{S}^{-1}\|_s^2 \|\Sigma_x - \widehat{S}\|_s}{1 - \|\widehat{S}^{-1}\|_s \|\Sigma_x - \widehat{S}\|_s}.$$

Moreover, since $\|\Sigma_x - \widehat{S}\|_s \leq d \|\Sigma_x - \widehat{S}\|_{\max}$, we obtain $\|\nabla_{\Theta} \mathbb{L}_T(\Theta_0)\|_{\max} \leq Ld\sqrt{\frac{\log(d^2)}{T}}$, for $L > 0$ sufficiently large. Hence, for a sample size $T \geq Md^2 \log(d^2) \max(R^2, k_0) \vee L \log(d^2)^{2/\gamma-1}$, then (10) holds with probability at least $1 - O(\exp(-\log(d^2))) - o(\exp(-\log(d^2)))$. \square

Proof of Corollary 3.5. First, let us consider the series expansion of the logarithm of a p -square positive definite matrix A , given as $\log(A) = -\sum_{k=1}^{\infty} \frac{1}{k} (I_p - A)^k$. Obviously, when the spectral radius of $I_p - A$ is strictly inferior to one, this expansion does exist. Moreover, the inverse of A can be developed as a Neumann series $A^{-1} = \sum_{k=0}^{\infty} (I_p - A)^k$. Let us consider the first order differential of $\mathbb{L}_{T,\eta}(\Theta)$:

$$\begin{aligned} d\mathbb{L}_{T,\eta}(\Theta) &= \eta \text{tr} \left((d\Theta/\nu) \log(\Theta/\nu) + (\Theta/\nu) (d \log(\Theta/\nu)) \right) \\ &\quad + (2\eta - 1) \text{tr} \left(- (d\Theta/\nu) \right) - \eta \text{tr} \left(\log(\widehat{S}^{-1}/\nu) (d\Theta/\nu) \right) - (1 - \eta) \text{tr} \left((d \log(\Theta/\nu)) (\widehat{S}^{-1}/\nu) \right) \end{aligned}$$

Now, let us treat the trace of $(\Theta/\nu) d \log(\Theta/\nu)$. Following Abadir and Magnus (2006) (see exercise 13.31),

we have

$$\begin{aligned} \text{tr} \left((\Theta/\nu) (d \log(\Theta/\nu)) \right) &= \text{tr} \left(\sum_{k=1}^{\infty} \frac{1}{k} \left\{ \sum_{l=1}^k (I_p - \Theta/\nu)^{l-1} (d\Theta/\nu) (I_p - \Theta/\nu)^{k-l} \right\} (\Theta/\nu) \right) \\ &= \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \text{tr} \left((I_p - \Theta/\nu)^{l-1} (d\Theta/\nu) (I_p - \Theta/\nu)^{k-l} (\Theta/\nu) \right) \\ &= \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \text{tr} \left((I_p - \Theta/\nu)^{k-l} (\Theta/\nu) (I_p - \Theta/\nu)^{l-1} (d\Theta/\nu) \right). \end{aligned}$$

In the same spirit,

$$\text{tr} \left((\widehat{S}^{-1}/\nu) (d \log(\Theta/\nu)) \right) = \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \text{tr} \left((I_p - \Theta/\nu)^{k-l} (\widehat{S}^{-1}/\nu) (I_p - \Theta/\nu)^{l-1} (d\Theta/\nu) \right).$$

Now let us consider the cases $\eta = 0, 1$ separately.

(i) $\eta = 1$: the differential becomes

$$d\mathbb{L}_{T,1}(\Theta) = \text{tr}\left((d\Theta/\nu) \log(\Theta/\nu) + (\Theta/\nu)(d \log(\Theta/\nu)) - (d\Theta/\nu) - \log(\widehat{S}^{-1}/\nu)(d\Theta/\nu)\right).$$

Thus, the differential becomes

$$d\mathbb{L}_{T,1}(\Theta) = \text{tr}\left([\log(\Theta/\nu) - \log(\widehat{S}^{-1}/\nu)](d\Theta/\nu) + \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k (I_p - \Theta/\nu)^{k-l} (\Theta/\nu) (I_p - \Theta/\nu)^{l-1} (d\Theta/\nu) - (d\Theta/\nu)\right).$$

Hence, by the trace properties and the von Neumann based series expansion of the inverse matrix, the gradient in vector and matrix forms are, respectively,

$$\nabla_{\theta} \mathbb{L}_{T,1}(\Theta) = \text{vec}(\log(\Theta/\nu) - \log(\widehat{S}^{-1}/\nu))/\nu, \quad \nabla_{\Theta} \mathbb{L}_{T,1}(\Theta) = \frac{1}{\nu} \left(\log(\Theta/\nu) - \log(\widehat{S}^{-1}/\nu) \right).$$

Let us now focus on the Hessian matrix, where we only need to focus on Starting from $\text{tr}\left((d\Theta/\nu) \log(\Theta/\nu)\right)/\nu$.

The second order differential becomes

$$d^2 \mathbb{L}_{T,1}(\Theta) = \text{tr}\left((d\Theta/\nu)(d \log(\Theta/\nu))\right)/\nu = \text{tr}\left((d\Theta/\nu) \left[\sum_{k=1}^{\infty} \frac{1}{k} \left\{ \sum_{l=1}^k (I_p - \Theta/\nu)^{l-1} (d\Theta/\nu) (I_p - \Theta/\nu)^{k-l} \right\} \right]\right).$$

Using the Hessian identification, we obtain

$$\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,\eta}(\Theta) = \sum_{k=1}^{\infty} \frac{1}{k} \left\{ \sum_{l=1}^k (I_p - \Theta/\nu)^{k-l} \otimes (I_p - \Theta/\nu)^{l-1} \right\} / \nu^2$$

By 10.20 (b) of Abadir and Magnus (2006), we have

$$f_T(\Theta) = \text{vec}(\Gamma)^\top \nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,\eta}(\Theta) \text{vec}(\Gamma) = \text{tr}\left(\sum_{k=1}^{\infty} \frac{1}{k} \left\{ \sum_{l=1}^k (I_p - \Theta/\nu)^{k-l} \Gamma (I_p - \Theta/\nu)^{l-1} \Gamma \right\}\right) / \nu^2.$$

where $\Theta = \Theta_0 + s\Gamma$, $\Gamma = \Theta_1 - \Theta_0$ with $\Theta_1 \in \Omega^{\text{vn}}$. Since Θ_0 and Θ_1 can be diagonalized in the same basis, idest $\Theta_0 = P\Lambda_0 P^{-1}$ and $\Theta_1 = P\Lambda_1 P^{-1}$, with P invertible containing the eigenvectors and $\Lambda_j = \text{diag}(\lambda_{1,j}, \dots, \lambda_{d,j})$, $j = 0, 1$, a diagonal matrix containing the eigenvalues of Λ_j , $j = 0, 1$. Hence, Θ can also be diagonalized in the same basis with eigenvalues $\tilde{\lambda}_k$, $k = 1, \dots, d$, so that we can manipulate the trace of diagonal matrices. Hence $\text{tr}\left((I_p - \Theta/\nu)^{k-l} \Gamma (I_p - \Theta/\nu)^{l-1} \Gamma\right) = \text{tr}\left((I_p - \Theta/\nu)^{k-1} \Gamma^2\right)$.

As a consequence,

$$\begin{aligned} f_T(\Theta) &\geq \sum_{k=1}^{\infty} \frac{1}{\nu^2} \sum_{n=1}^p (1 - \tilde{\lambda}_n/\nu)^{k-1} (\lambda_{n,1} - \lambda_{n,0})^2 \\ &\geq \sum_{n=1}^p \frac{(\lambda_{n,1} - \lambda_{n,0})^2}{\tilde{\lambda}_n \nu} \geq \sum_{n=1}^p \frac{(\lambda_{n,1} - \lambda_{n,0})^2}{\nu (\lambda_{\max}(\Theta_0) \vee \lambda_{\max}(\Theta_1))} \geq \|\Gamma\|_F^2 \frac{1}{\nu d}. \end{aligned}$$

We deduce

$$(\theta - \theta_0)^\top \nabla_{\theta^\top}^2 \mathbb{L}_{T,1}(\bar{\Theta})(\theta - \theta_0) \geq \|\text{vec}(\Theta) - \text{vec}(\Theta_0)\|_2^2 \frac{1}{\nu d}.$$

The RSC parameters are thus given by $\alpha_1 = \alpha_2 = 1/(\nu d)$, $\tau_1 = \tau_2 = 0$.

(ii) $\eta = 0$: the differential becomes

$$\begin{aligned} \mathbf{d}\mathbb{L}_{T,0}(\Theta) &= \text{tr}\left(\mathbf{d}\Theta/\nu\right) - \text{tr}\left(\mathbf{d}\log(\Theta/\nu)(\widehat{S}^{-1}/\nu)\right) \\ &= \text{tr}\left(\mathbf{d}\Theta/\nu\right) - \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \text{tr}\left((I_p - \Theta/\nu)^{k-l}(\widehat{S}^{-1}/\nu)(I_p - \Theta/\nu)^{l-1}(\mathbf{d}\Theta/\nu)\right). \end{aligned} \quad (29)$$

The gradient in matrix form becomes

$$\begin{aligned} \nabla_{\Theta} \mathbb{L}_{T,0}(\Theta) &= \text{tr}(I_p)/\nu - \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \text{tr}\left((I_p - \Theta/\nu)^{k-l}(\widehat{S}^{-1}/\nu)(I_p - \Theta/\nu)^{l-1}\right)/\nu \\ &= \text{tr}(I_p)/\nu - \sum_{k=1}^{\infty} \text{tr}\left((I_p - \Theta/\nu)^{k-1}(\widehat{S}^{-1}/\nu)\right)/\nu \\ &= \text{tr}(I_p)/\nu - \text{tr}\left((\Theta/\nu)^{-1}\widehat{S}^{-1}/\nu\right)/\nu. \end{aligned}$$

As for the Hessian, we aim at extracting the form $\text{tr}(L(\mathbf{d}\Theta)M(\mathbf{d}\Theta))$ for L (resp. M) any square $m \times m$ matrix (resp. $n \times n$). Applying the differential operator on (29), we obtain

$$\begin{aligned} \mathbf{d}^2 \mathbb{L}_{T,0}(\Theta) &= \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \left[\text{tr}\left(\sum_{j=1}^{k-l} (I_p - \Theta/\nu)^{j-1}(\mathbf{d}\Theta/\nu)(I_p - \Theta/\nu)^{k-l-j}(\widehat{S}^{-1}/\nu)(I_p - \Theta/\nu)^{l-1}(\mathbf{d}\Theta/\nu)\right) \right] \\ &+ \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \left[\text{tr}\left((I_p - \Theta/\nu)^{k-l}(\widehat{S}^{-1}/\nu) \sum_{i=1}^{l-1} (I_p - \Theta/\nu)^{i-1}(\mathbf{d}\Theta/\nu)(I_p - \Theta/\nu)^{l-1-i}(\mathbf{d}\Theta/\nu)\right) \right] \\ &= \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \sum_{j=1}^{k-l} \left[\text{tr}\left((I_p - \Theta/\nu)^{j-1}(\mathbf{d}\Theta/\nu)(I_p - \Theta/\nu)^{k-l-j}(\widehat{S}^{-1}/\nu)(I_p - \Theta/\nu)^{l-1}(\mathbf{d}\Theta/\nu)\right) \right] \\ &+ \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \sum_{i=1}^{l-1} \left[\text{tr}\left((I_p - \Theta/\nu)^{k-l}(\widehat{S}^{-1}/\nu)(I_p - \Theta/\nu)^{i-1}(\mathbf{d}\Theta/\nu)(I_p - \Theta/\nu)^{l-1-i}(\mathbf{d}\Theta/\nu)\right) \right]. \end{aligned}$$

By Hessian identification, we thus deduce

$$\begin{aligned} \nabla_{\theta^\top}^2 \mathbb{L}_{T,0}(\Theta) &= \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \sum_{j=1}^{k-l} \frac{1}{2} \left[\left\{ (I_p - \Theta/\nu)^{j-1} \otimes (I_p - \Theta/\nu)^{k-l-j}(\widehat{S}^{-1}/\nu)(I_p - \Theta/\nu)^{l-1} \right\} \right. \\ &\quad \left. + \left\{ (I_p - \Theta/\nu)^{j-1} \otimes (I_p - \Theta/\nu)^{l-1}(\widehat{S}^{-1}/\nu)(I_p - \Theta/\nu)^{k-l-j} \right\} \right] / \nu^2 \\ &+ \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \sum_{i=1}^{l-1} \frac{1}{2} \left[\left\{ (I_p - \Theta/\nu)^{k-l}(\widehat{S}^{-1}/\nu)(I_p - \Theta/\nu)^{i-1} \otimes (I_p - \Theta/\nu)^{l-1-i} \right\} \right. \\ &\quad \left. + \left\{ (I_p - \Theta/\nu)^{i-1}(\widehat{S}^{-1}/\nu)(I_p - \Theta/\nu)^{k-l} \otimes (I_p - \Theta/\nu)^{l-1-i} \right\} \right] / \nu^2. \end{aligned}$$

Consequently:

$$\begin{aligned}
f_T(\Theta) &= \text{vec}(\Gamma)^\top \nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta) \text{vec}(\Gamma) \\
&= \text{vec}(\Gamma)^\top \left(\sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \sum_{j=1}^{k-l} \frac{1}{2} \left[\left\{ (I_p - \Theta/\nu)^{j-1} \otimes (I_p - \Theta/\nu)^{k-l-j} (\widehat{S}^{-1}/\nu) (I_p - \Theta/\nu)^{l-1} \right\} \right. \right. \\
&\quad \left. \left. + \left\{ (I_p - \Theta/\nu)^{j-1} \otimes (I_p - \Theta/\nu)^{l-1} (\widehat{S}^{-1}/\nu) (I_p - \Theta/\nu)^{k-l-j} \right\} \right] / \nu^2 \right. \\
&\quad \left. + \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \sum_{i=1}^{l-1} \frac{1}{2} \left[\left\{ (I_p - \Theta/\nu)^{k-l} (\widehat{S}^{-1}/\nu) (I_p - \Theta/\nu)^{i-1} \otimes (I_p - \Theta/\nu)^{l-1-i} \right\} \right. \right. \\
&\quad \left. \left. + \left\{ (I_p - \Theta/\nu)^{i-1} (\widehat{S}^{-1}/\nu) (I_p - \Theta/\nu)^{k-l} \otimes (I_p - \Theta/\nu)^{l-1-i} \right\} \right] / \nu^2 \right) \text{vec}(\Gamma).
\end{aligned}$$

Using the properties of the minimum eigenvalue of the Kronecker product, of the product of positive definite matrices and convergence of geometric series, we obtain

$$\begin{aligned}
f_T(\Theta) &\geq \frac{1}{\nu^2} \|\Gamma\|_F^2 \left\{ \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \sum_{j=1}^{k-l} \lambda_{\min}(I_p - \Theta/\nu)^{k-2} \lambda_{\min}(\widehat{S}^{-1}/\nu) + \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \sum_{i=1}^{l-1} \lambda_{\min}(I_p - \Theta/\nu)^{k-2} \lambda_{\min}(\widehat{S}^{-1}/\nu) \right\} \\
&\geq \frac{1}{\nu^2} \|\Gamma\|_F^2 \lambda_{\min}(\widehat{S}^{-1}/\nu) \left\{ \sum_{k=1}^{\infty} (k-1) \lambda_{\min}(I_p - \Theta/\nu)^{k-2} \right\} \\
&\geq \frac{1}{\nu^2} \|\Gamma\|_F^2 \lambda_{\min}(\widehat{S}^{-1}/\nu) \left[1 / (1 - \lambda_{\min}(I_p - \Theta/\nu))^2 \right] \\
&\geq \|\Gamma\|_F^2 \lambda_{\min}(\widehat{S}^{-1}/\nu) \{ \lambda_{\max}(\Theta_0) + 1 \}^{-2},
\end{aligned}$$

since $1 - \lambda_{\min}(I_p - \Theta/\nu) \leq \lambda_{\max}(\Theta/\nu)$. We then deduce

$$(\theta - \theta_0)^\top \nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,\eta}(\bar{\Theta})(\theta - \theta_0) \geq \|\text{vec}(\Theta) - \text{vec}(\Theta_0)\|_2^2 \lambda_{\min}(\widehat{S}^{-1}/\nu) \{ \lambda_{\max}(\Theta_0) + 1 \}^{-2}.$$

The RSC parameters are $\alpha_1 = \alpha_2 = \lambda_{\min}(\widehat{S}^{-1}/\nu) \{ \lambda_{\max}(\Theta_0) + 1 \}^{-2}$, and $\tau_1 = \tau_2 = 0$.

Let us evaluate the probability of satisfying (14) for $\eta = 0$, where the gradient is given by

$$\nabla_{\Theta} \mathbb{L}_{T,0}(\Theta) = \left(I_d - (\widehat{S}^{-1}/\nu)(\Theta/\nu)^{-1} \right) / \nu = \left(I_d - \widehat{S}^{-1} \Theta^{-1} \right) / \nu.$$

Thus, we have

$$\|\nabla_{\Theta} \mathbb{L}_{T,0}(\Theta_0)\|_{\max} = \|(\Theta_0 - \widehat{S}^{-1}) \Theta_0^{-1} / \nu\|_{\max} \leq \|\Theta_0^{-1} / \nu\|_s \|\Theta_0 - \widehat{S}^{-1}\|_s.$$

Using the same steps as in the least squares case, we deduce $\|\Theta_0 - \widehat{S}^{-1}\|_{\max} \leq Ld\sqrt{\frac{\log(d^2)}{T}}$, with high probability and for $L > 0$ large enough. Then:

$$\|\nabla_{\Theta} \mathbb{L}_{T,0}(\Theta_0)\|_{\max} \leq Ld\sqrt{\frac{\log(d^2)}{T}} \|\Theta_0^{-1}/\nu\|_s,$$

with probability at least $1 - O(\exp(-\log(d^2))) - o(\exp(-\log(d^2)))$.

□

Proof of Corollary 3.6. We first establish the RSC condition. The gradient is defined as

$$\nabla_{\theta} \mathbb{L}_T(\Theta) = \frac{1}{2} \text{vec}(\Theta \widehat{S} + \widehat{S} \Theta - 2I_d), \quad \nabla_{\Theta} \mathbb{L}_T(\Theta) = \frac{1}{2} (\Theta \widehat{S} + \widehat{S} \Theta - 2I_d).$$

As for the Hessian, by identification, we have

$$\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta) = \frac{1}{2} (\widehat{S} \otimes I_d + I_d \otimes \widehat{S}).$$

For some $\Theta_1 \in \Omega$ and $u \in [0, 1]$, let $\Theta = \Theta_0 + u\Gamma$ with $\Gamma = \Theta_1 - \Theta_0$. Then

$$f_T(\Theta) := \text{vec}(\Gamma)^\top \nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta) \text{vec}(\Gamma) \geq \|\Gamma\|_F^2 \frac{1}{2} [\lambda_{\min}(\widehat{S} \otimes I_d + I_d \otimes \widehat{S})] \geq \|\Gamma\|_F^2 \lambda_{\min}(\widehat{S}).$$

For $\|\Gamma\|_F^2 = \|\Theta - \Theta_0\|_F^2$, we thus deduce

$$(\theta - \theta_0)^\top \nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\bar{\Theta})(\theta - \theta_0) \geq \lambda_{\min}(\widehat{S}) \|\theta - \theta_0\|_2^2,$$

for any $\bar{\Theta}$ located between Θ and Θ_0 . The RSC condition is then satisfied with coefficients $\alpha_1 = \lambda_{\min}(\widehat{S})$

and $\alpha_2 = \alpha_1$, $\tau_1 = \tau_2 = 0$.

We now evaluate the probability of satisfying (17). We have

$$\|\nabla_{\Theta} \mathbb{L}_T(\Theta_0)\|_{\max} = \left\| \frac{1}{2} \Theta_0 (\widehat{S} - \Theta_0^{-1}) + \frac{1}{2} (\widehat{S} - \Theta_0^{-1}) \Theta_0 \right\|_{\max} \leq \|\Theta_0\|_s \|\widehat{S} - \Sigma_x\|_s.$$

Hence, by Lemma 3.2, we have

$$\|\nabla_{\Theta} \mathbb{L}_T(\Theta_0)\|_{\max} \leq L \|\Theta_0\|_s d \sqrt{\frac{\log(d^2)}{T}},$$

for $L > 0$ large enough. Hence, for a sample size $T \geq M \|\Theta_0\|_s^2 d^2 \log(d^2) \max(R^2, k_0) \vee L \log(d^2)^{2/\gamma-1}$, then

(18) holds with probability at least $1 - O(\exp(-\log(d^2))) - o(\exp(-\log(d^2)))$

□

To establish the proofs for support recovery in subsection 3.3, the key point is to show that strict dual feasibility holds. To do so and for each sparse estimator case, following the PDW construction of Wainwright (2009) or Ravikumar et al. (2011), we construct a theoretical estimator optimized over the true subset (hence it is not possible to compute it empirically).

Proof of Corollary 3.7. The oracle estimator defined in (19) becomes for the Stein's loss:

$$\widehat{\Theta}^{\mathbf{g}, \mathcal{O}} := \arg \min_{\theta \in \mathbb{R}^{|\mathcal{A}|}: \Theta \in \Omega} \left\{ \mathbb{L}_{T,1}(\Theta) \right\} = \arg \min_{\theta \in \mathbb{R}^{|\mathcal{A}|}: \Theta \in \Omega} \left\{ \text{tr}(\widehat{\mathcal{S}}\Theta) - \log(|\Theta|) \right\},$$

Proof of point (i). We highlight that our approach differs from the proof methods of Theorem 1 of Ravikumar et al. (2011) or Corollary 4 of Loh and Wainwright (2017) dedicated to Gaussian loss based sparse precision matrix, which are based on the Brouwer's fixed point Theorem, in that our proof strategy shares the same spirit as in the proof of Corollaries 2 and 3 respectively dedicated to Linear regression with corrupted covariates and GLM of Loh and Wainwright (2017) to construct the estimator $\widehat{\Theta}_{\mathcal{A}}$ such that $\text{supp}(\widehat{\Theta}) \subseteq \mathcal{A}$ and $\widehat{\Theta}_{\mathcal{A}}$ is a zero subgradient point of

$$\widehat{\Theta}_{\mathcal{A}}^{\mathbf{g}} = \arg \min_{\Theta: \Theta \in \Omega, \text{supp}(\Theta) \subseteq \text{supp}(\Theta_0)} \left\{ \mathbb{L}_{T,1}(\Theta) + \mathbf{p}(\lambda_T, \theta) \right\},$$

where $\mathbb{L}_{T,1}(\cdot)$ and Ω are defined as in (5) for $\eta = 1$, and μ -amenable penalty functions (hence LASSO, SCAD and MCP). The latter amenability property is a key point since we require the incoherence condition in the μ -amenable case. We show that strict dual feasibility holds for such statistical problem. Now by the zero gradient condition (23) of the PDW step, we obtain

$$\nabla_{\theta} \mathbb{L}_{T,1}(\widehat{\Theta}) - \nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0) + \nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0) - \nabla \mathbf{q}(\lambda_T, \text{vec}(\widehat{\Theta}^{\mathbf{g}})) + \lambda_T \widehat{\mathbf{z}} = 0.$$

This implies

$$\widehat{\mathbf{K}} \text{vec}(\widehat{\Theta}^{\mathbf{g}} - \Theta_0) + \nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0) - \nabla \mathbf{q}(\lambda_T, \text{vec}(\widehat{\Theta}^{\mathbf{g}})) + \lambda_T \widehat{\mathbf{z}} = 0,$$

with $\widehat{\mathbf{K}} = \int_0^1 \nabla_{\theta\theta^{\top}}^2 \mathbb{L}_{T,1}(\Theta_0 + u(\widehat{\Theta}^{\mathbf{g}} - \Theta_0)) du$. Equivalently, we have

$$\begin{pmatrix} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}} & \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}^c} \\ \widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}} & \widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}^c} \end{pmatrix} \begin{pmatrix} \text{vec}(\widehat{\Theta}^{\mathbf{g}} - \Theta_0)_{\mathcal{A}} \\ \mathbf{0} \end{pmatrix} + \begin{pmatrix} \nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}} - \nabla \mathbf{q}(\lambda_T, \text{vec}(\widehat{\Theta}^{\mathbf{g}})_{\mathcal{A}}) \\ \nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}^c} - \nabla \mathbf{q}(\lambda_T, \text{vec}(\widehat{\Theta}^{\mathbf{g}})_{\mathcal{A}^c}) \end{pmatrix} + \lambda_T \begin{pmatrix} \widehat{\mathbf{z}}_{\mathcal{A}} \\ \widehat{\mathbf{z}}_{\mathcal{A}^c} \end{pmatrix} = \mathbf{0}.$$

Consequently, we obtain

$$\begin{aligned} \widehat{z}_{\mathcal{A}^c} &= \frac{1}{\lambda_T} \left\{ \nabla \mathbf{q}(\lambda_T, \text{vec}(\widehat{\Theta}^g))_{\mathcal{A}^c} - \nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}^c} + \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \left[\nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}} \right. \right. \\ &\quad \left. \left. - \nabla \mathbf{q}(\lambda_T, \text{vec}(\widehat{\Theta}^g))_{\mathcal{A}} + \lambda_T \widehat{z}_{\mathcal{A}} \right] \right\}. \end{aligned}$$

Using the regularity condition (v), $\nabla \mathbf{q}(\lambda_T, \text{vec}(\widehat{\Theta}^g))_{\mathcal{A}^c} = \nabla \mathbf{q}(\lambda_T, \mathbf{0}_{\mathcal{A}^c}) = \mathbf{0}_{\mathcal{A}^c}$. This implies

$$\widehat{z}_{\mathcal{A}^c} = \frac{1}{\lambda_T} \left\{ -\nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}^c} + \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \left[\nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}} - \nabla \mathbf{q}(\lambda_T, \text{vec}(\widehat{\Theta}^g))_{\mathcal{A}} + \lambda_T \widehat{z}_{\mathcal{A}} \right] \right\}.$$

Taking the ℓ_{∞} -norm, we obtain

$$\begin{aligned} \|\widehat{z}_{\mathcal{A}^c}\|_{\infty} &\leq \frac{1}{\lambda_T} \left\| -\nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}^c} + \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}} \right\|_{\infty} \\ &\quad + \frac{1}{\lambda_T} \left\| \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} (\lambda_T \widehat{z}_{\mathcal{A}} - \nabla_{\theta} \mathbf{q}(\lambda_T, \widehat{\theta}_{\Theta})_{\mathcal{A}}) \right\|_{\infty} \\ &\leq \frac{1}{\lambda_T} \left\| -\nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}^c} + \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}} \right\|_{\infty} + \|\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1}\|_{\infty}, \end{aligned}$$

where we used $\|\lambda_T \widehat{z}_{\mathcal{A}} - \nabla_{\theta} \mathbf{q}(\lambda_T, \text{vec}(\widehat{\Theta}^g))_{\mathcal{A}}\|_{\infty} = \|\nabla_{\theta} \mathbf{q}(\lambda_T, \text{vec}(\widehat{\Theta}^g))_{\mathcal{A}}\|_{\infty} \leq \lambda_T$ from Lemma 8 of Loh and Wainwright (2017). Furthermore, we have

$$\left\| -\nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}^c} + \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}} \right\|_{\infty} \leq \left\| \nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}^c} \right\|_{\infty} + \|\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}}\|_{\infty}.$$

First, let us consider $\|\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}}\|_{\infty}$. To do so, we first control the following quantity:

$$\begin{aligned} &\|\widehat{\mathbf{K}} - \nabla_{\theta\theta^{\top}}^2 \mathbb{L}_{T,1}(\Theta_0)\|_{\infty} \\ &= \left\| \int_0^1 (\nabla_{\theta\theta^{\top}}^2 \mathbb{L}_T(\Theta_0 + s(\widehat{\Theta}^g - \Theta_0)) - \nabla_{\theta\theta^{\top}}^2 \mathbb{L}_{T,1}(\Theta_0)) ds \right\|_{\infty} \\ &\leq \int_0^1 \left\| (\nabla_{\theta\theta^{\top}}^2 \mathbb{L}_T(\Theta_0 + s(\widehat{\Theta}^g - \Theta_0)) - \nabla_{\theta\theta^{\top}}^2 \mathbb{L}_{T,1}(\Theta_0)) \right\|_{\infty} ds \\ &= \int_0^1 \left\| \left\{ (\Theta_0 + s(\widehat{\Theta}^g - \Theta_0))^{-1} \otimes (\Theta_0 + s(\widehat{\Theta}^g - \Theta_0))^{-1} \right\} - \left\{ \Theta_0^{-1} \otimes \Theta_0^{-1} \right\} \right\|_{\infty} ds. \end{aligned}$$

Now for any $s \in [0, 1]$, using the consistency result from Corollary 3.3, we have

$$\left\| (\Theta_0 + s(\widehat{\Theta}^g - \Theta_0)) - \Theta_0 \right\|_{\infty} = s \|\widehat{\Theta}^g - \Theta_0\|_{\infty} \leq d \|\widehat{\Theta}^g - \Theta_0\|_{\max} \leq d \|\widehat{\Theta}^g - \Theta_0\|_F \leq Ld\sqrt{k_0} \sqrt{\frac{\log(d^2)}{T}},$$

for L sufficiently large. Hence, by the same reasoning as before on the norm on the inverse matrix difference:

$$\left\| (\Theta_0 + s(\widehat{\Theta}^g - \Theta_0))^{-1} - \Theta_0^{-1} \right\|_{\infty} \leq Ld\sqrt{k_0} \sqrt{\frac{\log(d^2)}{T}}.$$

Thus, using Lemma 13 of Loh and Wainwright (2017) on the upper bound of $\|A \otimes A - B \otimes B\|_{\infty}$, we obtain

$$\|\widehat{\mathbf{K}} - \nabla_{\theta\theta^{\top}}^2 \mathbb{L}_{T,1}(\Theta_0)\|_{\infty} \leq Ld\sqrt{k_0} \sqrt{\frac{\log(d^2)}{T}},$$

which provides the bound that hold with high probability with $L > 0$ when restricting to \mathcal{A} :

$$v_1 := \|\widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}} - \nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}^c\mathcal{A}}\|_\infty \leq L\sqrt{k_0^2 \frac{\log(d^2)}{T}}, \quad v_2 := \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - (\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0))_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty \leq L\sqrt{k_0^2 \frac{\log(d^2)}{T}},$$

by union bound and since $\|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}} - \nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}\mathcal{A}}\|_\infty \leq L\sqrt{k_0^2 \frac{\log(d^2)}{T}}$. Now we bound $\|\widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}}\|_\infty$

to prove the success of the PDW method. To do so, we consider the expansion

$$\|\widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}}\|_\infty \leq M_1 + M_2, \quad \text{where:}$$

$$\begin{aligned} M_1 &:= \|\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}^c\mathcal{A}} \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}}\|_\infty, \\ M_2 &:= \|\{\widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}^c\mathcal{A}} \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\} \nabla_\theta \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}}\|_\infty. \end{aligned}$$

First, let us highlight that $\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)] = \nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)$ since it involves no empirical estimator such as the sample variance covariance matrix. Then, we have

$$M_1 \leq \|\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}^c\mathcal{A}} \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty \|\nabla_\theta \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}}\|_\infty \leq \omega K \sqrt{\frac{\log(d^2)}{T}},$$

with high probability and using the incoherence condition $\|\mathbf{K}_{0,\mathcal{A}^c\mathcal{A}} \mathbf{K}_{0,\mathcal{A}\mathcal{A}}^{-1}\| \leq \omega$. As for M_2 , we have

$$\begin{aligned} M_2 &\leq \left\{ \left\| \left(\widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}^c\mathcal{A}} \right) \left(\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1} \right) \right\|_\infty \right. \\ &\quad + \left\| \left(\widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}^c\mathcal{A}} \right) \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1} \right\|_\infty \\ &\quad \left. + \left\| \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}^c\mathcal{A}} \left(\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1} \right) \right\|_\infty \right\} \|\nabla_\theta \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}}\|_\infty \\ &\leq \left\{ v_1 v_2 + v_1 \|\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty + v_2 \|\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}^c\mathcal{A}}\|_\infty \right\} \|\nabla_\theta \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}}\|_\infty \\ &\leq \left\{ v_1 v_2 + v_1 \|\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty + v_2 \|\Sigma_x\|_\infty^2 \right\} \|\nabla_\theta \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}}\|_\infty \\ &\leq K \sqrt{k_0^2 \|\Sigma_x\|_\infty^4 \frac{\log(d^2)}{T}} \|\nabla_\theta \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}}\|_\infty, \end{aligned}$$

with high probability and K sufficiently large, where we used $\|\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}^c\mathcal{A}}\|_\infty \leq \|\Sigma_x\|_\infty^2$. Then:

$$\|\widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}}\|_\infty \leq L \sqrt{k_0^2 \|\Sigma_x\|_\infty^4 \frac{\log(d^2)}{T}} \sqrt{\frac{\log(d^2)}{T}},$$

with high probability under the scaling behaviour $T > M k_0^2 \|\Sigma_x\|_\infty^4 \log(d^2)$ with M sufficiently large. More-

over, using the incoherence condition and

$$\begin{aligned}
& \|\widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}}\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty \\
& \leq \|\widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}}\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2\mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}^c\mathcal{A}}\mathbb{E}[\nabla_{\theta\theta^\top}^2\mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty + \|\mathbb{E}[\nabla_{\theta\theta^\top}^2\mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}^c\mathcal{A}}\mathbb{E}[\nabla_{\theta\theta^\top}^2\mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty \\
& \leq L\sqrt{k_0^2\|\Sigma_x\|_\infty^4\frac{\log(d^2)}{T}} + \omega.
\end{aligned}$$

Thus we have for $L > 0$ sufficiently large

$$\|\mathbf{z}_{\mathcal{A}^c}\|_\infty \leq \frac{1}{\lambda_T} \left(L_1\sqrt{k_0^3\|\Sigma_x\|_\infty^4\frac{\log(d^2)}{T}} + L_2\sqrt{\frac{\log(d^2)}{T}} \right) + L_3\sqrt{k_0^3\|\Sigma_x\|_\infty^4\frac{\log(d^2)}{T}} + \omega,$$

for $L_1, L_2, L_3 > 0$. Then strict dual feasibility of Theorem 5.1 is satisfied when $\frac{1}{1-\omega}L\sqrt{\frac{\log(d^2)}{T}} \leq \lambda_T$, under the scaling $T > Ck_0^2\|\Sigma_x\|_\infty^4\log(d^2)$.

We now focus on the ℓ_∞ bound to apply point (i) of Theorem 5.2. We have

$$\begin{aligned}
& \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1}\nabla_\theta\mathbb{L}_{T,1}(\Theta_0)_\mathcal{A}\|_\infty \\
& \leq \|\{\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2\mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\}\nabla_\theta\mathbb{L}_{T,1}(\Theta_0)_\mathcal{A}\|_\infty + \|\mathbb{E}[\nabla_{\theta\theta^\top}^2\mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\nabla_\theta\mathbb{L}_{T,1}(\Theta_0)_\mathcal{A}\|_\infty \\
& \leq \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2\mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty\|\nabla_\theta\mathbb{L}_{T,1}(\Theta_0)_\mathcal{A}\|_\infty + \|\mathbb{E}[\nabla_{\theta\theta^\top}^2\mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\nabla_\theta\mathbb{L}_{T,1}(\Theta_0)_\mathcal{A}\|_\infty \\
& \leq C_1\sqrt{k_0^3\frac{\log(d^2)}{T}}\sqrt{\frac{\log(d^2)}{T}} + C_2\sqrt{\|\Sigma_x\|_\infty^4\frac{\log(d^2)}{T}},
\end{aligned}$$

with $C_1, C_2 > 0$ and using our previous upper bounds. We proved

$$\begin{aligned}
\|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - (\mathbb{E}[\nabla_{\theta\theta^\top}^2\mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}\mathcal{A}})^{-1}\|_\infty & \leq \sqrt{k_0}\|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - (\mathbb{E}[\nabla_{\theta\theta^\top}^2\mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}\mathcal{A}})^{-1}\|_s \\
& \leq L\sqrt{k_0^2\frac{\log(d^2)}{T}} \leq \beta_\infty.
\end{aligned}$$

Hence

$$\|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty \leq \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - (\mathbb{E}[\nabla_{\theta\theta^\top}^2\mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}\mathcal{A}})^{-1}\|_\infty + \|(\mathbb{E}[\nabla_{\theta\theta^\top}^2\mathbb{L}_{T,1}(\Theta_0)]_{\mathcal{A}\mathcal{A}})^{-1}\|_\infty \leq 2\beta_\infty.$$

Consequently, by part (i) of Theorem 5.2, we obtain for $\tilde{L} > 0$

$$\|\widehat{\Theta}^g - \Theta_0\|_{\max} \leq \tilde{L}\sqrt{\frac{\log(d^2)}{T}} + \lambda_T\beta_\infty.$$

Proof of point (ii). The same approach as in the proof of (i) can be applied. Since the regularizer is assumed to be (μ, ζ) -amenable, we have by Lemma 5 of Loh and Wainwright (2017) that $\lambda_T\widehat{\mathbf{z}}_{\mathcal{A}} -$

$\nabla_{\theta} \mathbf{q}(\lambda_T, \text{vec}(\widehat{\Theta})_{\mathcal{A}}) = 0$. Hence we have

$$\|\widehat{\mathbf{z}}_{\mathcal{A}^c}\|_{\infty} \leq \frac{1}{\lambda_T} \left\| -\nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}^c} + \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}} \right\|_{\infty}.$$

Following the same steps as in the proof of part (i), by upper bounding $\|\nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}^c}\|_{\infty}$ and $\|\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \nabla_{\theta} \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}}\|_{\infty}$ we establish strict dual feasibility by Proposition 5.3. Then the remainder follows from part (ii) of Theorem 5.2. Point (ii) highlights the gain of non-convex penalties: the incoherence condition can be relaxed. \square

Proof of Corollary 3.8. We remind that the oracle estimator is defined as

$$\widehat{\Theta}^{\text{ls}, \mathcal{O}} = \arg \min_{\theta \in \mathbb{R}^{|\mathcal{A}|}: \Theta \in \Omega} \left\{ \mathbb{L}_T(\Theta) \right\} = \arg \min_{\theta \in \mathbb{R}^{|\mathcal{A}|}: \Theta \in \Omega} \left\{ \|\Theta - \widehat{S}^{-1}\|_F^2 \right\}.$$

We first establish strict dual feasibility for

$$\widehat{\Theta}_{\mathcal{A}}^{\text{ls}} = \arg \min_{\Theta: \Theta \in \Omega, \text{supp}(\Sigma) \subseteq \text{supp}(\Theta_0)} \left\{ \mathbb{L}_T(\Theta) + \mathbf{p}(\lambda_T, \theta) \right\},$$

where $\mathbb{L}_T(\cdot)$ and Ω are defined as in (8). We only focus on (μ, ζ) -amenable penalty functions: indeed, should we consider μ -amenable penalty functions, the incoherence condition is not satisfied. We have

$$\widehat{\mathbf{z}}_{\mathcal{A}^c} = \frac{1}{\lambda_T} \left\{ -\nabla_{\theta} \mathbb{L}_T(\Theta_0)_{\mathcal{A}^c} + \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \left[\nabla_{\theta} \mathbb{L}_T(\Theta_0)_{\mathcal{A}} - \nabla_{\theta} \mathbf{q}(\lambda_T, \text{vec}(\widehat{\Theta}^{\text{ls}})_{\mathcal{A}}) + \lambda_T \widehat{\mathbf{z}}_{\mathcal{A}} \right] \right\}.$$

Taking the ℓ_{∞} -norm, we obtain

$$\begin{aligned} \|\widehat{\mathbf{z}}_{\mathcal{A}^c}\|_{\infty} &\leq \frac{1}{\lambda_T} \left\| -\nabla_{\theta} \mathbb{L}_T(\Theta_0)_{\mathcal{A}^c} + \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \nabla_{\theta} \mathbb{L}_T(\Theta_0)_{\mathcal{A}} \right\|_{\infty} \\ &\quad + \frac{1}{\lambda_T} \left\| \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} (\lambda_T \widehat{\mathbf{z}}_{\mathcal{A}} - \nabla_{\theta} \mathbf{q}(\lambda_T, \text{vec}(\widehat{\Theta}^{\text{ls}})_{\mathcal{A}})) \right\|_{\infty} \\ &\leq \frac{1}{\lambda_T} \left\| -\nabla_{\theta} \mathbb{L}_T(\Theta_0)_{\mathcal{A}^c} + \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \nabla_{\theta} \mathbb{L}_T(\Theta_0)_{\mathcal{A}} \right\|_{\infty}, \end{aligned}$$

where we used $\lambda_T \widehat{\mathbf{z}}_{\mathcal{A}} - \nabla_{\theta} \mathbf{q}(\lambda_T, \text{vec}(\widehat{\Theta})_{\mathcal{A}}) = 0$ from Lemma 8 of Loh and Wainwright (2017). Then:

$$\|\widehat{\mathbf{z}}_{\mathcal{A}^c}\|_{\infty} \leq \frac{1}{\lambda_T} \left\| 2\text{vec}(\widehat{S}^{-1} - \Theta_0)_{\mathcal{A}^c} + \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} 2\text{vec}(\Theta_0 - \widehat{S}^{-1})_{\mathcal{A}} \right\|_{\infty}.$$

Since $\widehat{\mathbf{K}} = I_d$, we have with probability at least $1 - O(\exp(-\log(d^2))) - o(\exp(-\log(d^2)))$ that

$$\|\widehat{\mathbf{z}}_{\mathcal{A}^c}\|_{\infty} \leq \frac{2}{\lambda_T} \left(\|\text{vec}(\widehat{S}^{-1} - \Theta_0)_{\mathcal{A}^c}\|_{\infty} + \|\text{vec}(\Theta_0 - \widehat{S}^{-1})_{\mathcal{A}}\|_{\infty} \right) \leq \frac{L}{\lambda_T} \left(\sqrt{(d^2 - k_0) \frac{\log(d^2)}{T}} + \sqrt{k_0 \frac{\log(d^2)}{T}} \right),$$

using the arguments in the proof of Corollary 3.4 on the upper bound of the gradient. Provided $\lambda_T >$

$L \sqrt{\{(d^2 - k_0) \vee k_0\} \frac{\log(d^2)}{T}}$, strict dual feasibility holds and support recovery is satisfied by Theorem 5.1 of

Loh and Wainwright.

As for the ℓ_∞ -bound, using $\|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty = 1$, we have

$$\|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_T(\Theta)_{\mathcal{A}}\|_\infty = \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} 2\text{vec}(\Theta_0 - \widehat{S}^{-1})_{\mathcal{A}}\|_\infty \leq L \sqrt{k_0 \frac{\log(d^2)}{T}},$$

for $L > 0$ large enough, with probability $1 - O(\exp(-\log(d^2))) - o(\exp(-\log(d^2)))$. \square

Proof of Corollary 3.9. The oracle estimator for the von Neumann divergence is defined as

$$\widehat{\Theta}^{\text{vn}, \mathcal{O}} = \arg \min_{\theta \in \mathbb{R}^{|\mathcal{A}|}: \Theta \in \Omega^{\text{vn}}} \left\{ \mathbb{L}_{T,0}(\Theta) \right\} = \arg \min_{\theta \in \mathbb{R}^{|\mathcal{A}|}: \Theta \in \Omega^{\text{vn}}} \left\{ \text{tr}(\Theta - \log(\Theta/\nu) \widehat{S}^{-1})/\nu \right\}, \text{ with } \Omega^{\text{vn}} \text{ as in (11).}$$

Proof of point (i). The PDW construction is based on the estimator

$$\widehat{\Theta}_{\mathcal{A}}^{\text{vn}} = \arg \min_{\Theta: \Theta \in \Omega^{\text{vn}}, \text{supp}(\Sigma) \subseteq \text{supp}(\Theta_0)} \left\{ \mathbb{L}_{T,0}(\Theta) + \mathbf{p}(\lambda_T, \theta) \right\},$$

where $\mathbb{L}_{T,0}(\cdot)$ and Ω^{vn} are defined as in (11). The proof relies on the same steps as in the proof of Corollary 3.7: we need to show that strict dual feasibility holds for the previous criterion satisfied in the presence of μ -amenable penalty functions. Now based on the same steps as in point (i) of the proof of Corollary 3.7:

$$\begin{aligned} \|\widehat{\mathbf{z}}_{\mathcal{A}^c}\|_\infty &\leq \frac{1}{\lambda_T} \left\| -\nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}^c} + \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}} \right\|_\infty \\ &\quad + \frac{1}{\lambda_T} \left\| \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} (\lambda_T \widehat{\mathbf{z}}_{\mathcal{A}} - \nabla_\theta \mathbf{q}(\lambda_T, \text{vec}(\widehat{\Theta}^{\text{vn}})_{\mathcal{A}})) \right\|_\infty \\ &\leq \frac{1}{\lambda_T} \left\| -\nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}^c} + \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}} \right\|_\infty + \|\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty. \end{aligned}$$

Moreover, we have

$$\left\| -\nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}^c} + \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}} \right\|_\infty \leq \|\nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}^c}\|_\infty + \|\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}}\|_\infty.$$

First, let us consider $\|\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_{T,1}(\Theta_0)_{\mathcal{A}}\|_\infty$. To do so, we first control the following quantity:

$$\begin{aligned} &\widehat{\mathbf{K}} - \nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0) \\ &= \int_0^1 (\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0 + s(\widehat{\Theta}^{\text{vn}} - \Theta_0)) - \nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)) ds \\ &= \int_0^1 s \nabla_\theta (\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_s)) \text{vec}(\widehat{\Theta}^{\text{vn}} - \Theta_0) ds, \end{aligned}$$

where we used the mean value theorem with Θ_s lying between Θ_0 and $s(\widehat{\Theta} - \Theta_0)$. Such expansion is necessary due to the high non-linearity with respect to the parameters of the Hessian function. Let $w = d^2$,

for any $(\mathbf{u}, \mathbf{v}) \in \mathbb{R}^w \times \mathbb{R}^w$, we have

$$\begin{aligned}
& |\mathbf{u}^\top \{\widehat{\mathbf{K}} - \nabla_{\theta\theta}^2 \mathbb{L}_{T,0}(\Theta_0)\} \mathbf{v}| \\
&= \left| \int_0^1 \left\{ s \sum_{a,b,c=1}^w (\nabla_{abc}^3 \mathbb{L}_{T,0}(\Theta_s) \text{vec}(\widehat{\Theta}^{\text{vn}} - \Theta_0)_c \mathbf{u}_a \mathbf{v}_b) \right\} ds \right| \\
&\leq \int_0^1 s \left\{ \sum_{a,b,c=1}^w (\nabla_{abc}^3 \mathbb{L}_{T,0}(\Theta_s) \text{vec}(\widehat{\Theta}^{\text{vn}} - \Theta_0)_c \mathbf{u}_a \mathbf{v}_b) \right\} ds,
\end{aligned}$$

where ∇_{abc}^3 refers to the third order derivative with respect to θ . Let us now derive the third order derivative of $\mathbb{L}_{T,0}(\cdot)$. Element by element, for any $a, b, c = 1, \dots, w$, we have

$$\begin{aligned}
& \partial_{abc}^3 \mathbb{L}_{T,0}(\Theta) \\
&= - \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \sum_{j=1}^{k-l} \text{tr} \left(\sum_{x=1}^{j-1} (I_d - \Theta/\nu)^{x-1} (\partial_c \Theta/\nu) (I_d - \Theta/\nu)^{j-1-x} (\partial_b \Theta/\nu) (I_d - \Theta/\nu)^{k-l-j} (\widehat{S}^{-1}/\nu) (I_d - \Theta/\nu)^{l-1} (\partial_a \Theta/\nu) \right. \\
&\quad \left. + (I_d - \Theta/\nu)^{j-1} (\partial_b \Theta/\nu) \sum_{y=1}^{k-l-j} (I_d - \Theta/\nu)^{y-1} (\partial_c \Theta/\nu) (I_d - \Theta/\nu)^{k-l-j-y} (\widehat{S}^{-1}/\nu) (I_d - \Theta/\nu)^{l-1} (\partial_a \Theta/\nu) \right. \\
&\quad \left. + (I_d - \Theta/\nu)^{j-1} (\partial_b \Theta/\nu) (I_d - \Theta/\nu)^{k-l-j} (\widehat{S}^{-1}/\nu) \sum_{z=1}^{l-1} (I_d - \Theta/\nu)^{z-1} (\partial_c \Theta/\nu) (I_d - \Theta/\nu)^{l-1-z} (\partial_a \Theta/\nu) \right) \\
&\quad - \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \sum_{i=1}^{l-1} \text{tr} \left(\sum_{x=1}^{k-l} (I_d - \Theta/\nu)^{x-1} (\partial_c \Theta/\nu) (I_d - \Theta/\nu)^{k-l-x} (\widehat{S}^{-1}/\nu) (I_d - \Theta/\nu)^{i-1} (\partial_b \Theta/\nu) (I_d - \Theta/\nu)^{l-1-i} (\partial_a \Theta/\nu) \right. \\
&\quad \left. + (I_d - \Theta/\nu)^{k-l} (\widehat{S}^{-1}/\nu) \sum_{y=1}^{i-1} (I_d - \Theta/\nu)^{y-1} (\partial_c \Theta/\nu) (I_d - \Theta/\nu)^{i-1-y} (\partial_b \Theta/\nu) (I_d - \Theta/\nu)^{l-1-i} (\partial_a \Theta/\nu) \right. \\
&\quad \left. + (I_d - \Theta/\nu)^{k-l} (\widehat{S}^{-1}/\nu) (I_d - \Theta/\nu)^{i-1} (\partial_b \Theta/\nu) \sum_{z=1}^{l-1-i} (I_d - \Theta/\nu)^{z-1} (\partial_c \Theta/\nu) (I_d - \Theta/\nu)^{l-1-i-1} (\partial_a \Theta/\nu) \right) \\
&:= \sum_{r=1}^6 \Delta_{r,abc}.
\end{aligned}$$

Now by the Cauchy-Schwartz inequality, we have

$$|\nabla_{\theta} \left\{ \mathbf{u}^\top \nabla_{\theta\theta}^2 \mathbb{L}_{T,0}(\Theta_s) \mathbf{v} \right\} \text{vec}(\widehat{\Theta} - \Theta_0)|^2 \leq \left\{ \sum_{a,b,c=1}^w \partial_{abc}^3 \mathbb{L}_{T,0}(\Theta_s)^2 \mathbf{u}_a^2 \mathbf{v}_b^2 \right\} \|\widehat{\Theta}^{\text{vn}} - \Theta_0\|_F^2.$$

Regarding the third order derivative, taking the supremum on unit vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^{\text{card}(\mathcal{A})}$, using the Frobenius norm, we have the upper bound:

$$\sum_{a,b,c=1}^w \partial_{abc}^3 \mathbb{L}_{T,0}(\Theta_s)^2 \mathbf{u}_a^2 \mathbf{v}_b^2 \leq k_0 \left\{ \sqrt{k_0} \frac{1}{\nu^3} \|\widehat{S}^{-1}/\nu\|_F \sum_{k=1}^{\infty} \|I_d - \Theta_s/\nu\|_s^{k-3} (k^2 - 3k + 2) \right\}^2,$$

where we used $\|(I_d - \Theta_s/\nu)^l\|_F \leq \sqrt{k_0} \|(I_d - \Theta_s/\nu)^l\|_s \leq \sqrt{k_0} \|I_d - \Theta_s/\nu\|_s^l$ since we restrict to the \mathcal{A} block.

Now rewriting this series as a combination of geometric series, we deduce

$$\begin{aligned}
& \sum_{a,b,c=1}^w \partial_{abc}^3 \mathbb{L}_{T,0}(\Theta_s)^2 \mathbf{u}_a^2 \mathbf{v}_b^2 \\
& \leq k_0 \left\{ \frac{1}{\nu^3} \sqrt{k_0} \|\widehat{S}^{-1}/\nu\|_F \left[\|I_d - \Theta_s/\nu\|_s^{-1} \sum_{k=2}^{\infty} k(k-1) \|I_d - \Theta_s/\nu\|_s^{k-2} - 2 \|I_d - \Theta_s/\nu\|_s^{-2} \sum_{k=1}^{\infty} k \|I_d - \Theta_s/\nu\|_s^{k-1} \right. \right. \\
& \quad \left. \left. + 2 \|I_d - \Theta_s/\nu\|_s^{-2} \sum_{k=0}^{\infty} \|I_d - \Theta_s/\nu\|_s^k \right] \right\}^2 \\
& \leq k_0 \left\{ \frac{1}{\nu^3} \sqrt{k_0} \|\widehat{S}^{-1}/\nu\|_F \left[\|I_d - \Theta_s/\nu\|_s^{-1} \frac{2}{(1 - \|I_d - \Theta_s/\nu\|_s)^3} - 2 \|I_d - \Theta_s/\nu\|_s^{-2} \frac{1}{(1 - \|I_d - \Theta_s/\nu\|_s)^2} \right. \right. \\
& \quad \left. \left. + 2 \frac{1}{1 - \|I_d - \Theta_s/\nu\|_s} \right] \right\}^2 \\
& = 2k_0 \left\{ \frac{1}{\nu^3} \sqrt{k_0} \frac{\|\widehat{S}^{-1}/\nu\|_F}{1 - \|I_d - \Theta_s/\nu\|_s} \left[\|I_d - \Theta_s/\nu\|_s^{-1} \frac{1}{(1 - \|I_d - \Theta_s/\nu\|_s)^2} \right. \right. \\
& \quad \left. \left. - \|I_d - \Theta_s/\nu\|_s^{-2} \frac{2}{(1 - \|I_d - \Theta_s/\nu\|_s)} + 1 \right] \right\}^2 \\
& := k_0 \left\{ \frac{1}{\nu^3} \sqrt{k_0} \|\widehat{S}^{-1}/\nu\|_F \phi(\Theta_s, \nu) \right\}^2.
\end{aligned}$$

Now let us bound $1 - \|I_d - \Theta_s/\nu\|_s$. For ν large enough ensuring that the spectral norm of Θ/ν is smaller than 1, we have $\|I_d - \Theta_s/\nu\|_s = 1 - \lambda_{\min}(\Theta_s/\nu)$. The inequality $\frac{1}{1 - \|I_d - \Theta_s/\nu\|_s} \leq \lambda_{\min}(\Theta_s/\nu)^{-1}$ then holds.

Moreover,

$$\lambda_{\min}(\Theta_s/\nu) \geq \lambda_{\min}(\Theta_0/\nu) + s \lambda_{\min}((\widehat{\Theta}^{\text{vn}} - \Theta_0)/\nu) \geq \lambda_{\min}(\Theta_0/\nu) - \frac{s}{\nu} \|\widehat{\Theta}^{\text{vn}} - \Theta_0\|_F.$$

Thus using the upper bound with respect to $\|\widehat{\Theta}^{\text{vn}} - \Theta_0\|_F$ from Corollary 3.5, we obtain

$$\lambda_{\min}(\Theta_s/\nu) \geq \lambda_{\min}(\Theta_0/\nu) - \frac{s}{\nu} L \sqrt{k_0 \frac{\log(d^2)}{T}},$$

for $L > 0$ large enough. We deduce $\lambda_{\min}(\Theta_s/\nu) \geq \lambda_{\min}(\Theta_0/\nu) - \epsilon_{T,d} > 0$, with $\epsilon_{T,d} \rightarrow 0$ as $T, d \rightarrow \infty$ for a suitable scaling behaviour (T, d, k_0) . Consequently, we deduce $(1 - \|I_d - \Theta_s/\nu\|_s)^{-1} \leq (\lambda_{\min}(\Theta_0/\nu) - \epsilon_{T,d})^{-1}$.

Now let us treat $\|I_d - \Theta_s/\nu\|_s^{-1}$. We have

$$\|I_d - \Theta_s/\nu\|_s = 1 - \lambda_{\min}(\Theta_s/\nu) \geq 1 - \lambda_{\max}(\Theta_s/\nu) \geq 1 - \lambda_{\max}(\Theta_0/\nu) - \frac{1}{\nu} \|\widehat{\Theta}^{\text{vn}} - \Theta_0\|_F \geq 1 - \lambda_{\max}(\Theta_0/\nu) - \epsilon_{T,d},$$

with $\epsilon_{T,d} \rightarrow 0$ as $T, d \rightarrow \infty$ for a suitable scaling behaviour (T, d, k_0) . As a consequence, we obtain

$$\|I_d - \Theta_s/\nu\|_s^{-1} \leq \left(1 - \lambda_{\max}(\Theta_0/\nu) - \epsilon_{T,d}\right)^{-1}.$$

As for $\|\widehat{S}^{-1}\|_F$, we have $\|\widehat{S}^{-1}\|_F \leq \|\widehat{S}^{-1} - \Sigma_x^{-1}\|_F + \|\Sigma_x^{-1}\|_F$. Using the same arguments as in the proof of Corollary 3.4 for bounding $\|\widehat{S}^{-1} - \Sigma_x^{-1}\|_F$, we have $\|\widehat{S}^{-1} - \Sigma_x^{-1}\|_F \leq \sqrt{k_0}\|\widehat{S}^{-1} - \Sigma_x^{-1}\|_s \leq K\sqrt{k_0^2 \frac{\log(d^2)}{T}}$, with high probability and for $K > 0$ large enough. Putting all the pieces together, we deduce

$$\begin{aligned}
& \left\{ \sum_{a,b,c=1}^w \partial_{abc}^3 \mathbb{L}_{T,0}(\Theta_s)^2 \mathbf{u}_a^2 \mathbf{v}_b^2 \right\} \|\widehat{\Theta}^{\text{vn}} - \Theta_0\|_F^2 \\
& \leq Lk_0 \left[\frac{1}{\nu^4} \sqrt{k_0} \left(K \sqrt{k_0^2 \frac{\log(d^2)}{T}} + \|\Theta_0\|_F \right) \left(\lambda_{\min}(\Theta_0/\nu) - \epsilon_{T,d} \right)^{-1} \right. \\
& \quad \times \left\{ \left(1 - \lambda_{\max}(\Theta_0/\nu) - \epsilon_{T,d} \right)^{-1} \left(\lambda_{\min}(\Theta_0/\nu) - \epsilon_{T,d} \right)^{-2} - \left(1 - \lambda_{\max}(\Theta_0/\nu) - \epsilon_{T,d} \right)^{-2} \right. \\
& \quad \left. \left. \left(\lambda_{\min}(\Theta_0/\nu) - \epsilon_{T,d} \right)^{-1} + 1 \right\} \right]^2 k_0 \frac{\log(d^2)}{T} \\
& \leq Lk_0 \left[\frac{1}{\nu^4} \sqrt{k_0} \left(K \sqrt{k_0^2 \frac{\log(d^2)}{T}} + \|\Theta_0\|_F \right) \left(\lambda_{\min}(\Theta_0/\nu) - \epsilon_{T,d} \right)^{-2} \left(1 - \lambda_{\max}(\Theta_0/\nu) - \epsilon_{T,d} \right)^{-1} \right. \\
& \quad \times \left\{ \left(\lambda_{\min}(\Theta_0/\nu) - \epsilon_{T,d} \right)^{-1} - \left(1 - \lambda_{\max}(\Theta_0/\nu) - \epsilon_{T,d} \right)^{-1} \right. \\
& \quad \left. \left. + \left(\lambda_{\min}(\Theta_0/\nu) - \epsilon_{T,d} \right) \left(1 - \lambda_{\max}(\Theta_0/\nu) - \epsilon_{T,d} \right) \right\} \right]^2 k_0 \frac{\log(d^2)}{T}.
\end{aligned}$$

with high probability for $L > 0$ large enough. Now, let us treat $\epsilon_{T,d} \rightarrow 0$ under a suitable scaling (T, d) .

Using $\lambda_{\min}(\Theta_0/\nu)^{-1} = \nu \|\Sigma_x\|_s$, for $\tilde{L} > 0$ large enough, with high probability, the upper bound becomes

$$\begin{aligned}
& \left\{ \sum_{a,b,c=1}^w \partial_{abc}^3 \mathbb{L}_{T,0}(\Theta_s)^2 \mathbf{u}_a^2 \mathbf{v}_b^2 \right\} \|\widehat{\Theta}^{\text{vn}} - \Theta_0\|_F^2 \\
& \leq \tilde{L}k_0 \left[\frac{1}{\nu^4} \sqrt{k_0} \left(K \sqrt{k_0^2 \frac{\log(d^2)}{T}} + \|\Theta_0\|_F \right) \nu^3 \|\Sigma_x\|_s^2 \left(1 - \|\Theta_0/\nu\|_s \right)^{-1} \right]^2 k_0 \frac{\log(d^2)}{T}.
\end{aligned}$$

Thus, taking the supremum with respect to unit vector $\mathbf{u}, \mathbf{v} \in \mathbb{R}^A$, since $\|\Theta_0\|_F$ is of order k_0 , we obtain

$$\|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}} - \nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}\mathcal{A}}\|_s \leq C \sqrt{\frac{1}{\nu^2} k_0^7 \frac{\log(d^2)}{T}} \|\Sigma_x\|_s^4 \left(1 - \|\Theta_0/\nu\|_s \right)^{-2}. \quad (30)$$

for $C > 0$ sufficiently large. Furthermore, we have

$$\begin{aligned}
& \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}\|_s \\
& \leq \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}} - \nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}\mathcal{A}}\|_s + \|\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}\mathcal{A}} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}\|_s.
\end{aligned}$$

Let us control for $\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}\mathcal{A}} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}$. First, using the formula we obtained for

$\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta)$ in the proof of Corollary 3.5, the population level Hessian $\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]$ is

$$\begin{aligned}
& \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)] \\
&= \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \sum_{j=1}^{k-l} \frac{1}{2} \left[\left\{ (I_p - \Theta_0/\nu)^{j-1} \otimes (I_p - \Theta_0/\nu)^{k-l-j} \left(\mathbb{E}[\widehat{S}^{-1}] \right) / \nu (I_p - \Theta_0/\nu)^{l-1} \right\} \right. \\
&\quad \left. + \left\{ (I_p - \Theta_0/\nu)^{j-1} \otimes (I_p - \Theta_0/\nu)^{l-1} \left(\mathbb{E}[\widehat{S}^{-1}] \right) / \nu (I_p - \Theta_0/\nu)^{k-l-j} \right\} \right] / \nu^2 \\
&+ \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \sum_{i=1}^{l-1} \frac{1}{2} \left[\left\{ (I_p - \Theta_0/\nu)^{k-l} \left(\widehat{S}^{-1} - \mathbb{E}[\widehat{S}^{-1}] \right) / \nu (I_p - \Theta_0/\nu)^{i-1} \otimes (I_p - \Theta_0/\nu)^{l-1-i} \right\} \right. \\
&\quad \left. + \left\{ (I_p - \Theta_0/\nu)^{i-1} \left(\mathbb{E}[\widehat{S}^{-1}] \right) / \nu (I_p - \Theta_0/\nu)^{k-l} \otimes (I_p - \Theta_0/\nu)^{l-1-i} \right\} \right] / \nu^2. \tag{31}
\end{aligned}$$

Now by the properties on the Kronecker product, we obtain

$$\begin{aligned}
& \nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}\mathcal{A}} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}} \\
&= \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \sum_{j=1}^{k-l} \frac{1}{2} \left[\left\{ (I_p - \Theta_0/\nu)^{j-1} \otimes (I_p - \Theta_0/\nu)^{k-l-j} \left(\widehat{S}^{-1} - \mathbb{E}[\widehat{S}^{-1}] \right) / \nu (I_p - \Theta_0/\nu)^{l-1} \right\} \right. \\
&\quad \left. + \left\{ (I_p - \Theta_0/\nu)^{j-1} \otimes (I_p - \Theta_0/\nu)^{l-1} \left(\widehat{S}^{-1} - \mathbb{E}[\widehat{S}^{-1}] \right) / \nu (I_p - \Theta_0/\nu)^{k-l-j} \right\} \right] / \nu^2 \\
&+ \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \sum_{i=1}^{l-1} \frac{1}{2} \left[\left\{ (I_p - \Theta_0/\nu)^{k-l} \left(\widehat{S}^{-1} - \mathbb{E}[\widehat{S}^{-1}] \right) / \nu (I_p - \Theta_0/\nu)^{i-1} \otimes (I_p - \Theta_0/\nu)^{l-1-i} \right\} \right. \\
&\quad \left. + \left\{ (I_p - \Theta_0/\nu)^{i-1} \left(\widehat{S}^{-1} - \mathbb{E}[\widehat{S}^{-1}] \right) / \nu (I_p - \Theta_0/\nu)^{k-l} \otimes (I_p - \Theta_0/\nu)^{l-1-i} \right\} \right] / \nu^2.
\end{aligned}$$

First we consider $\mathbb{E}[\widehat{S}^{-1}]$. By a Taylor expansion, $\widehat{S}^{-1} = \Theta_0 + \Sigma_x^{-1} (\Sigma_x - \widehat{S}) \Sigma_x^{-1} + u_T$, where $u_T = o_p(1)$.

Thus, assuming $\mathbb{E}[\widehat{S}^{-1}] < \infty$, then $\mathbb{E}[\widehat{S}^{-1}] - \Theta_0 = o(1)$ for a sufficiently large T . Using the properties of the multiplicative norm of the Kronecker product, we obtain for T large enough

$$\begin{aligned}
& \|\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}\mathcal{A}} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}\|_s \\
&\leq \frac{1}{\nu^3} \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \sum_{j=1}^{k-l} \|I_p - \Theta_0/\nu\|_s^{k-2} \|\widehat{S}^{-1} - \Sigma_x^{-1}\|_s + \frac{1}{\nu^3} \sum_{k=1}^{\infty} \frac{1}{k} \sum_{l=1}^k \sum_{i=1}^{l-1} \|I_p - \Theta_0/\nu\|_s^{k-2} \|\widehat{S}^{-1} - \Sigma_x^{-1}\|_s \\
&\leq \frac{1}{\nu^3} \|\widehat{S}^{-1} - \Sigma_x^{-1}\|_s \sum_{k=1}^{\infty} (k-1) \|I_d - \Theta_0/\nu\|_s^{k-2}.
\end{aligned}$$

With high probability, we established $\|\widehat{S}^{-1} - \Sigma_x^{-1}\|_s \leq K \sqrt{k_0 \frac{\log(d^2)}{T}}$ when restricted to the \mathcal{A} block. As a consequence,

$$\begin{aligned}
\|\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}\mathcal{A}} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}\|_s &\leq \frac{1}{\nu^3} K \sqrt{k_0 \frac{\log(d^2)}{T}} \left(1 - \|I_d - \Theta_0/\nu\|_s\right)^{-2} \\
&\leq \frac{1}{\nu^3} K \sqrt{k_0 \frac{\log(d^2)}{T}} \lambda_{\min}(\Theta_0/\nu)^{-2} \\
&\leq K \sqrt{\frac{1}{\nu^2} \|\Sigma_x\|_s^4 k_0 \frac{\log(d^2)}{T}}.
\end{aligned}$$

for $L > 0$ sufficiently large since $\|I_d - \Theta_0/\nu\|_s < 1$. Thus, putting the pieces together, we deduce

$$\begin{aligned} & \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}\|_s \\ & \leq C \sqrt{\frac{1}{\nu^2} k_0^7 \frac{\log(d^2)}{T} \|\Sigma_x\|_s^4 \left(1 - \|\Theta_0/\nu\|_s\right)^{-2}} + K \sqrt{\frac{1}{\nu^2} \|\Sigma_x\|_s^4 k_0 \frac{\log(d^2)}{T}}. \end{aligned}$$

Thus we obtain

$$\|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}\|_s \leq L \sqrt{\frac{1}{\nu^2} k_0^7 \frac{\log(d^2)}{T} \|\Sigma_x\|_s^4 \left(1 - \|\Theta_0/\nu\|_s\right)^{-2}},$$

for L a sufficiently large constant. Now, using the same trick as in the proof of Corollary 3.4 for controlling for $\|\Theta_0^{-1} - \widehat{S}\|_s$, we obtain

$$\|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\|_s \leq L \sqrt{\frac{1}{\nu^2} k_0^7 \frac{\log(d^2)}{T} \|\Sigma_x\|_s^4 \left(1 - \|\Theta_0/\nu\|_s\right)^{-2}}. \quad (32)$$

Based on the same arguments for deriving (30), by union bound, we have

$$\max_{i \in \mathcal{A}^c} \|e_i^\top (\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} - \nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}^c \mathcal{A}})\|_2 \leq C \sqrt{\frac{1}{\nu^2} k_0^7 \frac{\log(d^2)}{T} \|\Sigma_x\|_s^4 \left(1 - \|\Theta_0/\nu\|_s\right)^{-2}},$$

which implies

$$\max_{i \in \mathcal{A}^c} \|e_i^\top (\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}^c \mathcal{A}})\|_2 \leq C \sqrt{\frac{1}{\nu^2} k_0^7 \frac{\log(d^2)}{T} \|\Sigma_x\|_s^4 \left(1 - \|\Theta_0/\nu\|_s\right)^{-2}}. \quad (33)$$

Now we are in a position to control the quantity:

$$\|\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_\mathcal{A}\|_\infty \leq M_1 + M_2, \quad \text{with}$$

$$\begin{aligned} M_1 & := \|\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}^c \mathcal{A}} \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}^c \mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_\mathcal{A}\|_\infty, \\ M_2 & := \left\| \left\{ \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}^c \mathcal{A}} \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1} \right\} \nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_\mathcal{A} \right\|_\infty. \end{aligned}$$

By assumption, the population level Hessian is bounded. As for the gradient, using Corollary 3.5, we obtain

with high probability that $M_1 \leq C \sqrt{k_0 \frac{\log(d^2)}{T}} \|\Theta_0^{-1}/\nu\|_s$. Regarding M_2 , we have

$$M_2 \leq \max_{i \in \mathcal{A}^c} \|e_i^\top \{ \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}^c \mathcal{A}} \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1} \} \|_2 \|\nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_\mathcal{A}\|_2. \quad (34)$$

Restricting to the true support, the following holds:

$$\begin{aligned} \|\nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_\mathcal{A}\|_2 & = \|\nabla_\Theta \mathbb{L}_{T,0}(\Theta_0)_\mathcal{A}\|_F \\ & \leq \sqrt{k_0} \|\nabla_\Theta \mathbb{L}_{T,0}(\Theta_0)_\mathcal{A}\|_s \leq \tilde{L} \sqrt{k_0} \sqrt{k_0 \frac{\log(d^2)}{T}} \|\Theta_0^{-1}/\nu\|_s = \tilde{L} \sqrt{k_0} \sqrt{\frac{1}{\nu^2} \|\Sigma_x\|_s^2 k_0 \frac{\log(d^2)}{T}}, \end{aligned} \quad (35)$$

for \tilde{L} a sufficiently large constant with high probability. Moreover

$$\begin{aligned} & \|e_i^\top \{\widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}}\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}^c\mathcal{A}} \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\}\|_2 \\ & \leq \|e_i^\top \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}^c\mathcal{A}} \Upsilon_1\|_2 + \|e_i^\top \Upsilon_2 \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\|_2 + \|e_i^\top \Upsilon_2 \Upsilon_1\|_2, \end{aligned} \quad (36)$$

with $\Upsilon_1 := \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}$, $\Upsilon_2 := \widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}^c\mathcal{A}}$. By inequalities (32) and (33), we obtain

$$\|\Upsilon_1\|_s \leq C \sqrt{\frac{1}{\nu^2} k_0^7 \frac{\log(d^2)}{T} \|\Sigma_x\|_s^4 (1 - \|\Theta_0/\nu\|_s)^{-2}}, \quad \max_{i \in \mathcal{A}^c} \|e_i^\top \Upsilon_2\|_2 \leq C \sqrt{\frac{1}{\nu^2} k_0^7 \frac{\log(d^2)}{T} \|\Sigma_x\|_s^4 (1 - \|\Theta_0/\nu\|_s)^{-2}}.$$

Hence, using inequalities (34), (35) and (36), we obtain for a sufficiently large constant \tilde{C} that

$$\begin{aligned} & \left\| \left\{ \widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}^c\mathcal{A}} \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1} \right\} \nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_\mathcal{A} \right\|_\infty \\ & \leq \tilde{C} \sqrt{k_0} \sqrt{\frac{1}{\nu^2} k_0 \frac{\log(d^2)}{T} \|\Sigma_x\|_s^2} \sqrt{\frac{1}{\nu^2} k_0^7 \frac{\log(d^2)}{T} \|\Sigma_x\|_s^4 (1 - \|\Theta_0/\nu\|_s)^{-2}}. \end{aligned}$$

Using the incoherence condition $\|\mathbf{K}_{0,\mathcal{A}^c\mathcal{A}} \mathbf{K}_{0,\mathcal{A}\mathcal{A}}^{-1}\|_\infty < \omega < 1$, we then obtain with high probability

$$\begin{aligned} & \|\widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_\mathcal{A}\|_\infty \\ & \leq \omega K_1 \sqrt{\frac{1}{\nu^2} \|\Sigma_x\|_s^2 k_0^2 \frac{\log(d^2)}{T}} + K_2 \sqrt{\frac{1}{\nu^4} k_0^9 \frac{\log(d^2)}{T} \|\Sigma_x\|_s^6 (1 - \|\Theta_0/\nu\|_s)^{-2}}, \end{aligned}$$

with $K_1, K_2 > 0$ large enough. Moreover, using the incoherence condition and

$$\begin{aligned} & \|\widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty \\ & \leq \sqrt{k_0} \|\widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}^c\mathcal{A}} \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty + \|\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}^c\mathcal{A}} \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty \\ & \leq \tilde{C} \sqrt{\frac{1}{\nu^4} k_0^9 \frac{\log(d^2)}{T} \|\Sigma_x\|_s^6 (1 - \|\Theta_0/\nu\|_s)^{-2}} + \omega. \end{aligned}$$

Thus, putting the pieces together, we have for $L_1, L_2, L_3 > 0$ sufficiently large

$$\begin{aligned} & \|\mathbf{z}_{\mathcal{A}^c}\|_\infty \\ & \leq \frac{1}{\lambda_T} \left(L_1 \sqrt{\frac{1}{\nu^4} k_0^9 \frac{\log(d^2)}{T} \|\Sigma_x\|_s^6 (1 - \|\Theta_0/\nu\|_s)^{-2}} + L_2 \sqrt{\frac{1}{\nu^2} (d^2 - k_0) \|\Sigma_x\|_s^2 \frac{\log(d^2)}{T}} \right) \\ & \quad + L_3 \sqrt{\frac{1}{\nu^2} k_0^7 \frac{\log(d^2)}{T} \|\Sigma_x\|_s^4 (1 - \|\Theta_0/\nu\|_s)^{-2}} + \omega. \end{aligned}$$

Hence, strict dual feasibility of Theorem 5.1 is satisfied when

$$\frac{1}{1-\omega} L \sqrt{\frac{\log(d^2)}{T}} \leq \lambda_T,$$

under the scaling $T > C \frac{1}{\nu^2} \left[(d^2 - k_0) \vee \frac{k_0^9}{\nu^2} \|\Sigma_x\|_s^4 \left(1 - \|\Theta_0/\nu\|_s\right)^{-2} \right] \|\Sigma_x\|_s^2 \log(d^2)$.

We now focus on the ℓ_∞ bound. We have

$$\begin{aligned} & \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}}\|_\infty \\ & \leq \|\{\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\} \nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}}\|_\infty + \|\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}}\|_\infty \\ & \leq \sqrt{k_0} \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\|_s \|\nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}}\|_\infty + \|\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}}\|_\infty \\ & \leq C_1 \sqrt{k_0} \sqrt{\frac{1}{\nu^2} k_0^7 \frac{\log(d^2)}{T} \|\Sigma_x\|_s^4 \left(1 - \|\Theta_0/\nu\|_s\right)^{-2}} \sqrt{\frac{1}{\nu^2} \|\Sigma_x\|_s^2 k_0 \frac{\log(d^2)}{T}} \\ & \quad + C_2 \sqrt{\frac{1}{\nu^2} \|\Sigma_x\|_s^2 k_0 \frac{\log(d^2)}{T}}, \end{aligned}$$

with $C_1, C_2 > 0$ using inequality (32). We proved

$$\begin{aligned} & \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - (\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}})^{-1}\|_\infty \\ & \leq \sqrt{k_0} \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - (\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}})^{-1}\|_s \leq \sqrt{\frac{1}{\nu^2} k_0^7 \frac{\log(d^2)}{T} \|\Sigma_x\|_s^4 \left(1 - \|\Theta_0/\nu\|_s\right)^{-2}} \leq \beta_\infty. \end{aligned}$$

Hence

$$\|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty \leq \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - (\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}})^{-1}\|_\infty + \|(\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_{T,0}(\Theta_0)]_{\mathcal{A}\mathcal{A}})^{-1}\|_\infty \leq 2\beta_\infty.$$

Consequently, by part (i) of Theorem 5.2, we obtain for $\tilde{L} > 0$

$$\|\widehat{\Theta}^{\text{vn}} - \Theta_0\|_{\max} \leq \tilde{L} \sqrt{\frac{\log(d^2)}{T}} + \lambda_T \beta_\infty.$$

Proof of point (ii). We follow the same proof as in point (ii) of Corollary 3.7. Hence we have

$$\|\widehat{\mathbf{z}}_{\mathcal{A}^c}\|_\infty \leq \frac{1}{\lambda_T} \|\nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}^c}\|_\infty + \|\widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}}\|_\infty.$$

Following the same steps as in the proof of part (i), by upper bounding $\|\nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}^c}\|_\infty$ and $\|\widehat{\mathbf{K}}_{\mathcal{A}^c\mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_{T,0}(\Theta_0)_{\mathcal{A}}\|_\infty$

we establish strict dual feasibility by Proposition 5.3. Then the remainder follows from part (ii) of Theorem

5.2.

□

Proof of Corollary 3.10. The oracle estimator is defined as

$$\widehat{\Theta}^{\text{dt}, \mathcal{O}} = \arg \min_{\theta \in \mathbb{R}^{|\mathcal{A}|}: \Theta \in \Omega} \left\{ \mathbb{L}_T(\Theta) \right\} = \arg \min_{\theta \in \mathbb{R}^{|\mathcal{A}|}: \Theta \in \Omega} \left\{ \frac{1}{2} \text{tr}(\Theta^2 \widehat{S} - \Theta) \right\}.$$

Proof of point (i). Strict dual feasibility is checked for

$$\widehat{\Theta}^{\text{dt}} = \arg \min_{\Theta: \Theta \in \Omega, \text{supp}(\Theta) \subseteq \text{supp}(\Theta_0)} \left\{ \mathbb{L}_T(\Theta) + \mathbf{p}(\lambda_T, \theta) \right\},$$

where $\mathbb{L}_T(\cdot)$ and Ω are defined in (16). Following the same steps as in the Stein's case in point (i), we aim at bounding

$$\|\widehat{z}_{\mathcal{A}^c}\|_\infty \leq \frac{1}{\lambda_T} \left\| -\nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}^c} + \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}} \right\|_\infty + \|\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1}\|_\infty,$$

with $\widehat{\mathbf{K}} = \int_0^1 \nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0 + u(\widehat{\Theta}^{\text{dt}} - \Theta_0)) du$. Now we have

$$\left\| -\nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}^c} + \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}} \right\|_\infty \leq \left\| \nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}^c} \right\|_\infty + \left\| \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}} \right\|_\infty.$$

Let us highlight that the Hessian matrix does not depend on the parameter Θ . Proceeding as in the proof of point (i) of Corollary 3.7, we obtain

$$\left\| \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}} \right\|_\infty \leq M_1 + M_2, \quad \text{where}$$

$$\begin{aligned} M_1 &= \left\| \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}^c \mathcal{A}} \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A} \mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}} \right\|_\infty, \\ M_2 &= \left\| \left\{ \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}^c \mathcal{A}} \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A} \mathcal{A}}^{-1} \right\} \nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}} \right\|_\infty. \end{aligned}$$

Note that $\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)] = \frac{1}{2} (\Sigma_x \otimes I_d + I_d \otimes \Sigma_x)$. Then, we have

$$M_1 \leq \left\| \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}^c \mathcal{A}} \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A} \mathcal{A}}^{-1} \right\|_\infty \left\| \nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}} \right\|_\infty \leq \omega K \sqrt{\|\Theta_0\|_s^2 k_0 \frac{\log(d^2)}{T}},$$

with high probability and using the incoherence condition, a condition identical to Zhang and Zou (2014),

Theorem 2. As for M_2 , we first consider:

$$\left\| \widehat{\mathbf{K}} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)] \right\|_\infty = \frac{1}{2} \left\| \left(\widehat{S} - \Sigma_x \otimes I_d \right) + \left(I_d \otimes \widehat{S} - \Sigma_x \right) \right\|_\infty \leq \left\| \widehat{S} - \Sigma_x \right\|_\infty \leq d \|\widehat{S} - \Sigma_x\|_{\max} \leq L \sqrt{d^2 \frac{\log(d^2)}{T}},$$

with high probability for $L > 0$ large enough. Then let us consider the probability:

$$\begin{aligned} & \forall \epsilon > 0, \mathbb{P}(\|\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}^c \mathcal{A}}\|_\infty > \epsilon) \\ & \leq \mathbb{P}(\max_{j \in \mathcal{A}^c} \sum_{l=1}^{k_0} |\widehat{\mathbf{K}}_{jl} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{jl}| > \epsilon) = \mathbb{P}(\max_{j \in \mathcal{A}^c} \sum_{l=1}^{k_0} |\widehat{S}_{jl} - \Sigma_{x,jl}| > \epsilon) \leq (d^2 - k_0) \sum_{l=1}^{k_0} \mathbb{P}(|\widehat{S}_{jl} - \Sigma_{x,jl}| > \epsilon/k_0). \end{aligned}$$

Then we deduce for $L > 0$ large enough and with a suitable sample size that

$$v_1 = \|\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}^c \mathcal{A}}\|_\infty, \quad v_2 = \|\widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A} \mathcal{A}}^{-1}\|_\infty,$$

are bounded by $L\sqrt{k_0 \frac{\log(d^2)}{T}}$ with high probability. Then

$$\begin{aligned} M_2 & \leq \left\{ \left\| \left(\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}^c \mathcal{A}} \right) \left(\widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A} \mathcal{A}}^{-1} \right) \right\|_\infty \right. \\ & \quad + \left\| \left(\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}^c \mathcal{A}} \right) \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A} \mathcal{A}}^{-1} \right\|_\infty \\ & \quad \left. + \left\| \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}^c \mathcal{A}} \left(\widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A} \mathcal{A}}^{-1} \right) \right\|_\infty \right\} \|\nabla_{\Theta} \mathbb{L}_T(\Theta_0)_{\mathcal{A}}\|_{\max} \\ & \leq \left\{ v_1 v_2 + v_1 \|\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A} \mathcal{A}}^{-1}\|_\infty + v_2 \|\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}^c \mathcal{A}}\|_\infty \right\} \|\nabla_{\Theta} \mathbb{L}_T(\Theta_0)_{\mathcal{A}}\|_{\max} \\ & \leq \left\{ v_1 v_2 + v_1 \|\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A} \mathcal{A}}^{-1}\|_\infty + v_2 \|\Sigma_x\|_\infty \right\} \|\nabla_{\Theta} \mathbb{L}_T(\Theta_0)_{\mathcal{A}}\|_{\max} \\ & \leq K \sqrt{k_0 \|\Sigma_x\|_\infty^2 \frac{\log(d^2)}{T}} \|\nabla_{\Theta} \mathbb{L}_T(\Theta_0)_{\mathcal{A}}\|_{\max}, \end{aligned}$$

with high probability and for K sufficiently large, where we used $\|\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}^c \mathcal{A}}\|_\infty \leq \|\Sigma_x\|_\infty$. Since

$\|\nabla_{\Theta} \mathbb{L}_T(\Theta_0)_{\mathcal{A}}\|_{\max} \leq K \sqrt{k_0 \|\Theta_0\|_s^2 \frac{\log(d^2)}{T}}$ with high probability, we conclude that

$$\|\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} \nabla_{\Theta} \mathbb{L}_T(\Theta_0)_{\mathcal{A}}\|_\infty \leq L \sqrt{k_0 \|\Sigma_x\|_\infty^2 \frac{\log(d^2)}{T}} \sqrt{\|\Theta_0\|_s^2 \frac{k_0 \log(d^2)}{T}}.$$

with high probability under the scaling $T > M k_0^2 \left\{ \|\Sigma_x\|_\infty^2 \vee \|\Theta_0\|_s^2 \right\} \log(d^2)$. Furthermore, by the incoherence condition

$$\begin{aligned} & \|\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1}\|_\infty \\ & \leq \left\| \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A} \mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}^c \mathcal{A}} \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A} \mathcal{A}}^{-1} \right\|_\infty + \left\| \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}^c \mathcal{A}} \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A} \mathcal{A}}^{-1} \right\|_\infty \\ & \leq L \sqrt{k_0 \|\Sigma_x\|_\infty^2 \frac{\log(d^2)}{T}} + \omega. \end{aligned}$$

Thus we have for $L > 0$ sufficiently large

$$\|\mathbf{z}_{\mathcal{A}^c}\|_\infty \leq \frac{1}{\lambda_T} \left(L_1 \sqrt{k_0^2 \left\{ \|\Sigma_x\|_\infty^2 \vee \|\Theta_0\|_s^2 \right\} \frac{\log(d^2)}{T}} + L_2 \sqrt{(d^2 - k_0) \|\Theta_0\|_s^2 \frac{\log(d^2)}{T}} \right) + L_3 \sqrt{k_0 \|\Sigma_x\|_\infty^2 \frac{\log(d^2)}{T}} + \omega,$$

for $L_1, L_2, L_3 > 0$. Hence, then strict dual feasibility of Theorem 5.1 is satisfied when

$$\frac{1}{1-\omega} L \sqrt{\frac{\log(d^2)}{T}} \leq \lambda_T,$$

under the scaling $T > C \max \left\{ k_0^2 \left[\|\Sigma_x\|_\infty^2 \vee \|\Theta_0\|_s^2 \right], (d^2 - k_0) \|\Theta_0\|_s^2 \right\} \log(d^2)$.

Let us consider the ℓ_∞ bound to apply point (i) of Theorem 5.2. We have

$$\begin{aligned} & \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}}\|_\infty \\ & \leq \|\{\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\} \nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}}\|_\infty + \|\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}}\|_\infty \\ & \leq \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - \mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty \|\nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}}\|_\infty + \|\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}}\|_\infty \\ & \leq C_1 \sqrt{k_0 \frac{\log(d^2)}{T}} \sqrt{k_0 \|\Theta_0\|_s^2 \frac{\log(d^2)}{T}} + C_2 \sqrt{k_0 \|\Theta_0\|_\infty^2 \frac{\log(d^2)}{T}}, \end{aligned}$$

with $C_1, C_2 > 0$ and using our previous upper bounds. We proved

$$\|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - (\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}\mathcal{A}})^{-1}\|_\infty \leq L \sqrt{k_0 \frac{\log(d^2)}{T}} \leq \beta_\infty.$$

Hence

$$\|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1}\|_\infty \leq \|\widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} - (\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}\mathcal{A}})^{-1}\|_\infty + \|(\mathbb{E}[\nabla_{\theta\theta^\top}^2 \mathbb{L}_T(\Theta_0)]_{\mathcal{A}\mathcal{A}})^{-1}\|_\infty \leq 2\beta_\infty.$$

Consequently, by part (i) of Theorem 5.2, we obtain

$$\|\widehat{\Theta}^g - \Theta_0\|_{\max} \leq \tilde{L} \sqrt{\frac{\log(d^2)}{T}} + \lambda_T \beta_\infty,$$

for $\tilde{L} > 0$.

Proof of point (ii). The same approach as in the proof of (i) can be applied. Since the regularizer is assumed to be (μ, ζ) -amenable, we have by Lemma 5 of Loh and Wainwright (2017) that $\lambda_T \widehat{\mathbf{z}}_{\mathcal{A}} - \nabla_\theta \mathbf{q}(\lambda_T, \text{vec}(\widehat{\Theta})_{\mathcal{A}}) = 0$. Hence we have

$$\|\widehat{\mathbf{z}}_{\mathcal{A}^c}\|_\infty \leq \frac{1}{\lambda_T} \|\nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}^c} + \widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}}\|_\infty.$$

Following the same steps as in the proof of part (i), by upper bounding $\|\nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}^c}\|_\infty$ and $\|\widehat{\mathbf{K}}_{\mathcal{A}^c \mathcal{A}} \widehat{\mathbf{K}}_{\mathcal{A}\mathcal{A}}^{-1} \nabla_\theta \mathbb{L}_T(\Theta_0)_{\mathcal{A}}\|_\infty$, we establish strict dual feasibility by Proposition 5.3. Then the remainder follows from part (ii) of Theorem

5.2. □