

Noise covariance estimation in multi-task high-dimensional linear models

KAI TAN ^{*1,a}, GABRIEL ROMON ^{2,c} and PIERRE C BELLEC ^{1,b}

¹Department of Statistics, Rutgers University, Piscataway, NJ 08854, USA, ^akai.tan@rutgers.edu,

^bpierre.bellec@rutgers.edu

²CREST, ENSAE, IP Paris, Palaiseau 91120 Cedex, France, ^cgabriel.romon@ensae.fr

This paper studies the multi-task high-dimensional linear regression models where the noise among different tasks is correlated, in the moderately high dimensional regime where sample size n and dimension p are of the same order. Our goal is to estimate the covariance matrix of the noise random vectors, or equivalently the correlation of the noise variables on any pair of two tasks. Treating the regression coefficients as a nuisance parameter, we leverage the multi-task elastic-net and multi-task lasso estimators to estimate the nuisance. By precisely understanding the bias of the squared residual matrix and by correcting this bias, we develop a novel estimator of the noise covariance that converges in Frobenius norm at the rate $n^{-1/2}$ when the covariates are Gaussian distributed with a known covariance matrix. This novel estimator is efficiently computable. Under suitable conditions, the proposed estimator of the noise covariance attains the same rate of convergence as the “oracle” estimator that knows in advance the regression coefficients of the multi-task model. The Frobenius error bounds obtained in this paper also illustrate the advantage of this new estimator compared to a method-of-moments estimator that does not attempt to estimate the nuisance. As byproducts of our techniques, we obtain estimates of the generalization error and out-of-sample error of the multi-task elastic-net and multi-task lasso estimators. Extensive simulation studies are carried out to illustrate the numerical performance of the proposed method.

Keywords: elastic-net; high dimensional analysis; lasso; multi-task model; noise covariance

1. Introduction

1.1. Model and estimation target

Consider a multi-task (also known as multi-response) linear model with T tasks and n i.i.d. observations $(\mathbf{x}_i, Y_{i1}, Y_{i2}, \dots, Y_{iT}), \forall i = 1, \dots, n$, where $\mathbf{x}_i \in \mathbb{R}^p$ is a random feature vector and Y_{i1}, \dots, Y_{iT} are responses in the model

$$\begin{aligned} Y_{it} &= \mathbf{x}_i^\top \boldsymbol{\beta}^{(t)} + E_{it} \quad \text{for each } t = 1, \dots, T; i = 1, \dots, n \quad (\text{scalar form}), \\ \mathbf{y}^{(t)} &= \mathbf{X} \boldsymbol{\beta}^{(t)} + \boldsymbol{\varepsilon}^{(t)} \quad \text{for each } t = 1, \dots, T \quad (\text{vector form}), \\ \mathbf{Y} &= \mathbf{X} \mathbf{B}^* + \mathbf{E} \quad (\text{matrix form}), \end{aligned} \tag{1}$$

where $\mathbf{X} \in \mathbb{R}^{n \times p}$ is the design matrix with rows $(\mathbf{x}_i^\top)_{i=1, \dots, n}$, $\mathbf{y}^{(t)} = (Y_{1t}, \dots, Y_{nt})^\top$ is the response vector for task t , $\boldsymbol{\varepsilon}^{(t)} = (E_{1t}, \dots, E_{nt})^\top$ is the noise vector for task t , $\boldsymbol{\beta}^{(t)} \in \mathbb{R}^p$ is an unknown fixed coefficient vector for task t . In matrix form, $\mathbf{Y} \in \mathbb{R}^{n \times T}$ is the response matrix with columns $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(T)}$, $\mathbf{E} \in \mathbb{R}^{n \times T}$ has columns $\boldsymbol{\varepsilon}^{(1)}, \dots, \boldsymbol{\varepsilon}^{(T)}$, and $\mathbf{B}^* \in \mathbb{R}^{p \times T}$ is an unknown coefficient matrix with columns $\boldsymbol{\beta}^{(1)}, \dots, \boldsymbol{\beta}^{(T)}$. The three forms in (1) are equivalent.

While the n vectors $(\mathbf{x}_i^\top, y_i^{(1)}, \dots, y_i^{(T)})_{i=1, \dots, n}$ of dimension $p + T$ are i.i.d., we assume that for each observation $i = 1, \dots, n$, the noise random variables E_{i1}, \dots, E_{iT} are centered and correlated. The

focus of the present paper is on estimation of the noise covariance matrix $\mathbf{S} \in \mathbb{R}^{T \times T}$, which has entries $S_{tt'} = \mathbb{E}[\varepsilon_1^{(t)} \varepsilon_1^{(t')}]$ for any pair $t, t' = 1, \dots, T$, or equivalently $\mathbf{S} = \frac{1}{n} \mathbb{E}[\mathbf{E}^\top \mathbf{E}]$.

The noise covariance plays a crucial role in multi-task linear models because it characterizes the noise level and correlation between different tasks: if tasks $t = 1, \dots, T$ represent time this captures temporal correlation; if tasks $t = 1, \dots, T$ represent different activation areas in the brain (e.g., [Bertrand et al. \(2019\)](#)) this captures spatial correlation.

Since \mathbf{S} is the estimation target, we view \mathbf{B}^* as an unknown nuisance parameter. If $\mathbf{B}^* = \mathbf{0}$, then $\mathbf{Y} = \mathbf{E}$, hence \mathbf{E} is directly observed and a natural estimator is the sample covariance $\frac{1}{n} \mathbf{E}^\top \mathbf{E}$. There are other possible choices for the sample covariance; ours coincides with the maximum likelihood estimator of the centered Gaussian model where the n samples are i.i.d. from $\mathcal{N}_T(\mathbf{0}, \mathbf{S})$. In the presence of a nuisance parameter $\mathbf{B}^* \neq \mathbf{0}$, the above sample covariance is not computable since we only observe (\mathbf{X}, \mathbf{Y}) and do not have access to \mathbf{E} . Thus we will refer to $\frac{1}{n} \mathbf{E}^\top \mathbf{E} \in \mathbb{R}^{T \times T}$ as the *oracle estimator* for \mathbf{S} , and its error $\frac{1}{n} \mathbf{E}^\top \mathbf{E} - \mathbf{S}$ will serve as a benchmark.

The nuisance parameter \mathbf{B}^* is not of interest by itself, but if an estimator $\hat{\mathbf{B}}$ is available that provides good estimation of \mathbf{B}^* , we would hope to leverage $\hat{\mathbf{B}}$ to estimate the nuisance and improve estimation of \mathbf{S} . For instance given an estimate $\hat{\mathbf{B}}$ such that $\|\mathbf{X}(\hat{\mathbf{B}} - \mathbf{B}^*)\|_{\text{F}}^2/n \rightarrow 0$, one may use the estimator

$$\hat{\mathbf{S}}_{(\text{naive})} = \frac{1}{n} (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})^\top (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}}) \quad (2)$$

to consistently estimate \mathbf{S} in Frobenius norm. We refer to this estimator as the *naive estimator* since it is obtained by simply replacing the noise \mathbf{E} in the oracle estimator $\frac{1}{n} \mathbf{E}^\top \mathbf{E}$ with the residual matrix $\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}}$. However, in the regime $p/n \rightarrow \gamma$ of interest in the present paper, the convergence $\|\mathbf{X}(\hat{\mathbf{B}} - \mathbf{B}^*)\|_{\text{F}}^2/n \rightarrow 0$ does not hold even in the case where $T = 1$ and where $\hat{\mathbf{B}}$ is chosen as the Ridge ([Dobriban and Wager, 2018](#)) or the Lasso ([Bayati and Montanari, 2012](#), [Miolane and Montanari, 2021](#)): The theory developed in these papers shows that $\|\mathbf{X}(\hat{\mathbf{B}} - \mathbf{B}^*)\|_{\text{F}}^2/n$ has a non-vanishing limit as p/n converges to a constant. Simulations in Section 4 will show that the naive estimator (2) presents a major bias for estimation of \mathbf{S} . One goal of this paper is to develop an estimator $\hat{\mathbf{S}}$ of \mathbf{S} by exploiting a commonly used estimator $\hat{\mathbf{B}}$ of the nuisance so that in the regime $p/n \rightarrow \gamma$ the error $\hat{\mathbf{S}} - \mathbf{S}$ is comparable to the benchmark $\frac{1}{n} \mathbf{E}^\top \mathbf{E} - \mathbf{S}$.

1.2. Related literature

If $T = 1$, the above model (1) reduces to the standard linear model with $\mathbf{X} \in \mathbb{R}^{n \times p}$ and response vector $\mathbf{y}^{(1)} \in \mathbb{R}^n$. We will refer to the $T = 1$ case as the single-task linear model and drop the superscript ⁽¹⁾ for brevity, i.e., $y_i = \mathbf{x}_i^\top \boldsymbol{\beta}^* + \varepsilon_i$, where ε_i are i.i.d. with mean 0, and unknown variance σ^2 . The coefficient vector $\boldsymbol{\beta}^*$ is typically assumed to be s -sparse, i.e., $\boldsymbol{\beta}^*$ has at most s nonzero entries. In this single-task linear model, estimation of noise covariance \mathbf{S} reduces to estimation of the noise variance $\sigma^2 = \mathbb{E}[\varepsilon_i^2]$, which has been studied in the literature. [Fan, Guo and Hao \(2012\)](#) proposed a consistent estimator for σ^2 based on a refitted cross validation method, which assumes the support of $\boldsymbol{\beta}^*$ is correctly recovered; [Belloni, Chernozhukov and Wang \(2011\)](#) and [Sun and Zhang \(2012\)](#) introduced square-root Lasso (scaled Lasso) to jointly estimate the coefficient $\boldsymbol{\beta}^*$ and noise variance σ^2 by

$$(\hat{\boldsymbol{\beta}}, \hat{\sigma}) = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p, \sigma > 0} \frac{\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2}{2n\sigma} + \frac{\sigma}{2} + \lambda_0 \|\boldsymbol{\beta}\|_1. \quad (3)$$

This estimator $\hat{\sigma}$ is consistent only when the prediction error $\|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\|^2/n$ goes to 0, which requires $s \log(p)/n \rightarrow 0$. Estimation of σ^2 without assumption on \mathbf{X} was proposed in [Yu and Bien \(2019\)](#) by

utilizing natural parameterization of the penalized likelihood of the linear model. Their estimator can be expressed as the minimizer of the Lasso problem: $\hat{\sigma}_\lambda^2 = \min_{\beta \in \mathbb{R}^p} \frac{1}{n} \|\mathbf{y} - \mathbf{X}\beta\|^2 + 2\lambda \|\beta\|_1$. Consistency of these estimators [Belloni, Chernozhukov and Wang \(2011, 2014\)](#), [Sun and Zhang \(2012\)](#), [Yu and Bien \(2019\)](#) requires $s \log(p)/n \rightarrow 0$ and does not hold in the high-dimensional proportional regime $p/n \rightarrow \gamma \in (0, \infty)$. For this proportional regime $p/n \rightarrow \gamma \in (0, \infty)$, under the assumption that \mathbf{x}_i are i.i.d. $\mathcal{N}(\mathbf{0}, \Sigma)$, [Dicker \(2014\)](#) introduced a method-of-moments estimator $\hat{\sigma}^2$ of σ^2 ,

$$\hat{\sigma}^2 = \frac{n+p+1}{n(n+1)} \|\mathbf{y}\|^2 - \frac{1}{n(n+1)} \|\Sigma^{-\frac{1}{2}} \mathbf{X}^\top \mathbf{y}\|^2, \quad (4)$$

which is unbiased, consistent, and asymptotically normal in high-dimensional linear models with Gaussian predictors and errors. Moreover, [Janson, Foygel Barber and Candès \(2017\)](#) developed an EigenPrism procedure for the same task as well as confidence intervals for σ^2 . The estimation procedures in these two papers don't attempt to estimate the nuisance parameter β^* , and require no sparsity on β^* and isometry structure, but assume $\|\Sigma^{\frac{1}{2}} \beta^*\|^2$ is bounded. Maximum Likelihood Estimators (MLEs) were studied in [Dicker and Erdogdu \(2016\)](#) for joint estimation of noise level and signal strength in high-dimensional linear models with fixed effects; they showed that a classical MLE for random-effects models may also be used effectively in fixed-effects models.

In the proportional regime, [Bayati, Erdogdu and Montanari \(2013\)](#), [Miolane and Montanari \(2021\)](#) used the Lasso to estimate the nuisance β^* and produce estimator for σ^2 . Their approach requires an uncorrelated Gaussian design assumption with $\Sigma = \mathbf{I}_p$. [Bellec \(2020\)](#) provided consistent estimators of a similar nature for σ^2 using more general M-estimators with convex penalty without requiring $\Sigma = \mathbf{I}_p$. In the special case of the squared loss, this estimator has the form [Bayati, Erdogdu and Montanari \(2013\)](#), [Bellec \(2020\)](#), [Miolane and Montanari \(2021\)](#)

$$\hat{\sigma}^2 = (n - \hat{\text{df}})^{-2} \{ \|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2 (n + p - 2\hat{\text{df}}) - \|\Sigma^{-\frac{1}{2}} (\mathbf{y} - \mathbf{X}\hat{\beta})\|^2 \}, \quad (5)$$

where $\hat{\text{df}} = \text{Tr}[(\partial/\partial \mathbf{y}) \mathbf{X}\hat{\beta}]$ denotes the degrees of freedom. This estimator coincides with the method-of-moments estimator in [Dicker \(2014\)](#) when $\hat{\beta} = \mathbf{0}$.

For multi-task high-dimensional linear model (1) with $T \geq 2$, the estimation of \mathbf{B}^* is studied in [Lounici et al. \(2011\)](#), [Obozinski, Wainwright and Jordan \(2011\)](#), [Simon, Friedman and Hastie \(2013\)](#). These works suggest to use a joint convex optimization problem over the tasks to estimate \mathbf{B}^* . A popular choice is the multi-task elastic-net, which solves the convex optimization problem

$$\hat{\mathbf{B}} = \arg \min_{\mathbf{B} \in \mathbb{R}^{p \times T}} \left(\frac{1}{2n} \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_{\text{F}}^2 + \lambda \|\mathbf{B}\|_{2,1} + \frac{\tau}{2} \|\mathbf{B}\|_{\text{F}}^2 \right), \quad (6)$$

where $\|\mathbf{B}\|_{2,1} = \sum_{j=1}^p \|\mathbf{B}^\top \mathbf{e}_j\|_2$, and $\|\cdot\|_{\text{F}}$ denotes the Frobenius norm of a matrix. This optimization problem can be efficiently solved by existing statistical packages, for instance, scikit-learn ([Pedregosa et al., 2011](#)), and glmnet ([Friedman, Hastie and Tibshirani, 2010](#)). Note that (6) is also referred to as multi-task (group) Lasso and multi-task Ridge if $\tau = 0$ and $\lambda = 0$, respectively. [Geer and Stucky \(2016\)](#) extended square-root Lasso ([Belloni, Chernozhukov and Wang, 2011](#)) and scaled Lasso ([Sun and Zhang, 2012](#)) to multi-task setting by solving the following problem

$$(\hat{\mathbf{B}}, \hat{\mathbf{S}}) = \arg \min_{\mathbf{B}, \mathbf{S} \succ 0} \left\{ \frac{1}{n} \text{Tr}((\mathbf{Y} - \mathbf{X}\mathbf{B})\mathbf{S}^{-\frac{1}{2}}(\mathbf{Y} - \mathbf{X}\mathbf{B})^\top) + \text{Tr}(\mathbf{S}^{\frac{1}{2}}) + 2\lambda_0 \|\mathbf{B}\|_1 \right\}, \quad (7)$$

where $\|\mathbf{B}\|_1 = \sum_{j,t} |B_{jt}|$. Note that the covariance estimator in (7) is constrained to be positive definite. [Molstad \(2022\)](#) studied the same problem and proposed to estimate \mathbf{S} by (2) with $\hat{\mathbf{B}}$ in (7), which is

consistent under Frobenius norm loss when $\|X(\widehat{\mathbf{B}} - \mathbf{B}^*)\|_F^2/n \rightarrow 0$. [Chen and Banerjee \(2017\)](#) considered another multi-response model with same coefficient for different tasks but different covariates. This literature investigated a similar joint estimation of coefficient and noise covariance. In a recent paper, [Bellec and Romon \(2021\)](#) studied the multi-task Lasso problem and proposed confidence intervals for single entries of \mathbf{B}^* and confidence ellipsoids for single rows of \mathbf{B}^* under the assumption that \mathbf{S} is proportional to the identity, which may be restrictive in practice. This literature generalizes degrees of freedom adjustments from single-task to multi-task models, which we will illustrate in Section 2.

Noise covariance estimation in the high dimensional multi-task linear model is a difficult problem. If the estimand \mathbf{S} is known to be diagonal, estimating \mathbf{S} reduces to the estimation of noise variance for each task, in which the existing methods for single-task high-dimensional linear models can be applied. Nonetheless, for general positive semi-definite matrix \mathbf{S} , the noise among different tasks may be correlated, hence the existing methods are not readily applicable, and a more careful analysis is called for to incorporate the correlation between different tasks. [Fourdrinier, Haddouche and Mezoued \(2021\)](#) considered estimating \mathbf{S} for the multi-task model (1) where rows of \mathbf{E} have elliptically symmetric distribution and in the classical regime $p \leq n$. However, their estimator has no statistical guarantee under Frobenius norm loss.

Recently, for the proportional regime $p/n \rightarrow \gamma \in (0, \infty)$, [Celentano and Montanari \(2021\)](#) generalized the estimator $\widehat{\sigma}^2$ in [Bayati, Erdogdu and Montanari \(2013\)](#) to the multi-task setting with $T = 2$. Their work covers correlated Gaussian designs, where a Lasso or Ridge regression is used to estimate $\boldsymbol{\beta}^{(1)}$ for the first task, and another Lasso or Ridge regression is used to estimate $\boldsymbol{\beta}^{(2)}$ for the second task. In other words, they estimate the coefficient vector for each task separately instead of using a multi-task estimator like (6). It is not trivial to adapt their estimator from the setting $T = 2$ to larger T , and allow T to increase with n . This present paper takes a different route and aims to fill this gap by proposing a novel noise covariance estimator with theoretical guarantees. Of course, our method applies directly to the 2-task linear model considered in [Celentano and Montanari \(2021\)](#).

1.3. Main contributions

The present paper introduces a novel estimator $\widehat{\mathbf{S}}$ in (11) of the noise covariance \mathbf{S} when the predictor covariance $\boldsymbol{\Sigma}$ is known. The proposed estimator $\widehat{\mathbf{S}}$ is shown to be a consistent estimator of \mathbf{S} under Frobenius norm, in the regime where p and n are of the same order. The estimator $\widehat{\mathbf{S}}$ is based on the multi-task elastic-net estimator $\widehat{\mathbf{B}}$ in (6) of the nuisance, and can be seen as a de-biased version of the naive estimator (2). The naive estimator (2) suffers from a strong bias in the regime where p and n are of the same order, and the estimator $\widehat{\mathbf{S}}$ is constructed by precisely understanding this bias and correcting it.

After introducing this novel estimator $\widehat{\mathbf{S}}$ in Definition 2.2 below, we prove several rates of convergence for the Frobenius error $\|\widehat{\mathbf{S}} - \mathbf{S}\|_F$, which is comparable, in terms of rate of convergence, to the benchmark $\|\frac{1}{n}\mathbf{E}^\top \mathbf{E} - \mathbf{S}\|_F$ under suitable assumptions.

As a by-product of the techniques developed for the construction of $\widehat{\mathbf{S}}$, we obtain estimates of the generalization error and out-of-sample error of $\widehat{\mathbf{B}}$, which are of independent interest and can be used for parameter tuning.

1.4. Notation

Basic notation and definitions that will be used in the rest of the paper are given here. Let $[n] = \{1, 2, \dots, n\}$ for all $n \in \mathbb{N}$. The vectors $\mathbf{e}_i \in \mathbb{R}^n, \mathbf{e}_j \in \mathbb{R}^p, \mathbf{e}_t \in \mathbb{R}^T$ denote the canonical basis vector of the corresponding index. We consider restrictions of vectors (resp., of matrices) by zeroing the

corresponding entries (resp., columns). More precisely, for $\mathbf{v} \in \mathbb{R}^p$ and index set $B \subset [p]$, $\mathbf{v}_B \in \mathbb{R}^p$ is the vector with $(\mathbf{v}_B)_j = 0$ if $j \notin B$ and $(\mathbf{v}_B)_j = v_j$ if $j \in B$. If $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $B \subset [p]$, $\mathbf{X}_B \in \mathbb{R}^{n \times p}$ is such that $(\mathbf{X}_B)\mathbf{e}_j = \mathbf{0}$ if $j \notin B$ and $(\mathbf{X}_B)\mathbf{e}_j = \mathbf{X}\mathbf{e}_j$ if $j \in B$. For a real vector $\mathbf{a} \in \mathbb{R}^p$, $\|\mathbf{a}\|$ denotes its Euclidean norm. For any matrix \mathbf{A} , \mathbf{A}^\dagger is its Moore–Penrose inverse; $\|\mathbf{A}\|_F$, $\|\mathbf{A}\|_{\text{op}}$, $\|\mathbf{A}\|_*$ denote its Frobenius, operator and nuclear norm, respectively. Let $\|\mathbf{A}\|_0$ be the number of non-zero rows of \mathbf{A} . Let $\mathbf{A} \otimes \mathbf{B}$ be the Kronecker product of \mathbf{A} and \mathbf{B} , and $\langle \mathbf{A}, \mathbf{B} \rangle = \text{Tr}(\mathbf{A}^\top \mathbf{B})$ is the Frobenius inner product for matrices of identical size. For \mathbf{A} symmetric, $\phi_{\min}(\mathbf{A})$ and $\phi_{\max}(\mathbf{A})$ denote its smallest and largest eigenvalues, respectively. Let \mathbf{I}_n denote the identity matrix of size n for all $n \in \mathbb{N}$. For a random sequence ξ_n , we write $\xi_n = O_P(a_n)$ if ξ_n/a_n is stochastically bounded. C denotes an absolute constant and $C(\tau, \gamma)$ stands for a generic positive constant depending on τ, γ ; their expression may vary from place to place.

1.5. Organization

The rest of the paper is organized as follows. Section 2 introduces our proposed estimator for noise covariance. Section 3 presents our main theoretical results on proposed estimator and some relevant estimators. Section 4 demonstrates through numerical experiments that our estimator outperforms several existing methods in the literature, which corroborates our theoretical findings in Section 3. Section 5 provides discussion and points out some future research directions. The appendix provides some technical results and proofs of all theoretical results in the paper. These technical results are further proved in the supplementary material [Tan, Romon and Bellec \(2023\)](#).

2. Estimating noise covariance, with possibly diverging number of tasks T

Before we can define our noise covariance estimator, we need to introduce the following building blocks. Let $\hat{\mathcal{J}} = \{k \in [p] : \hat{\mathbf{B}}^\top \mathbf{e}_k \neq \mathbf{0}\}$ denote the set of nonzero rows of $\hat{\mathbf{B}}$ in (6), and let $|\hat{\mathcal{J}}|$ denote the cardinality of $\hat{\mathcal{J}}$. For each $k \in \hat{\mathcal{J}}$, define $\mathbf{H}^{(k)} = \lambda \|\hat{\mathbf{B}}^\top \mathbf{e}_k\|^{-1} (\mathbf{I}_T - \hat{\mathbf{B}}^\top \mathbf{e}_k \mathbf{e}_k^\top \hat{\mathbf{B}} \|\hat{\mathbf{B}}^\top \mathbf{e}_k\|^{-2})$, which is the Hessian of the map $\mathbf{u} \mapsto \lambda \|\mathbf{u}\|$ at $\mathbf{u} = \hat{\mathbf{B}}^\top \mathbf{e}_k$ when $\mathbf{u} \neq \mathbf{0}$. Define $\mathbf{M} \in \mathbb{R}^{pT \times pT}$ by

$$\mathbf{M} = \mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + \tau n \mathbf{P}_{\hat{\mathcal{J}}}) + n \sum_{k \in \hat{\mathcal{J}}} (\mathbf{H}^{(k)} \otimes \mathbf{e}_k \mathbf{e}_k^\top), \quad (8)$$

where $\mathbf{P}_{\hat{\mathcal{J}}} = \sum_{k \in \hat{\mathcal{J}}} \mathbf{e}_k \mathbf{e}_k^\top \in \mathbb{R}^{p \times p}$. While the expression of \mathbf{M} looks a bit long, it is obtained by differentiating the multi-task estimate $\hat{\mathbf{B}}$ in (6) w.r.t. noise E_{it} , in the sense that $\mathbf{M} \frac{\partial \text{vec}(\hat{\mathbf{B}})}{\partial E_{it}} = (\mathbf{e}_t \otimes \mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{e}_i)$. The detailed derivation is presented in Proof of Lemma E.6. in supplementary file [Tan, Romon and Bellec \(2023\)](#). Define the residual matrix \mathbf{F} , the error matrix \mathbf{H} by

$$\mathbf{F} = \mathbf{Y} - \mathbf{X}\hat{\mathbf{B}}, \quad \mathbf{H} = \Sigma^{1/2}(\hat{\mathbf{B}} - \mathbf{B}). \quad (9)$$

To construct our estimator we also make use of the so-called interaction matrix $\hat{\mathbf{A}} \in \mathbb{R}^{T \times T}$.

Definition 2.1 ([Bellec and Romon \(2021\)](#)). The *interaction matrix* $\hat{\mathbf{A}} \in \mathbb{R}^{T \times T}$ of the estimator $\hat{\mathbf{B}}$ in (6) is defined by

$$\hat{\mathbf{A}} = \sum_{i=1}^n (\mathbf{I}_T \otimes \mathbf{e}_i^\top \mathbf{X}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \mathbf{X}^\top \mathbf{e}_i). \quad (10)$$

The matrix $\widehat{\mathbf{A}}$ was introduced in [Bellec and Romon \(2021\)](#), where it is used alongside the multi-task Lasso estimator ($\tau = 0$ in (6)). The interaction matrix $\widehat{\mathbf{A}}$ is essentially the derivatives of the mapping $\mathbf{E} \mapsto \mathbf{X}\widehat{\mathbf{B}}$, in the sense that $\widehat{\mathbf{A}}_{tt'} = \text{Tr}(\frac{\partial \mathbf{X}\widehat{\mathbf{B}}_{\mathbf{e}_t}}{\partial \mathbf{E}_{\mathbf{e}_{t'}}})$. Lemma E.6 in the supplementary material [Tan, Romon and Bellec \(2023\)](#) gives a formal statement of this explanation. It generalizes the degrees of freedom from [Stein \(1981\)](#) to the multi-task case. Intuitively, it captures the correlation between the residuals on different tasks ([Bellec and Romon, 2021](#), Lemma F.1). Our definition of the noise covariance estimator involves $\widehat{\mathbf{A}}$, although our statistical purposes differ greatly from the confidence intervals developed in [Bellec and Romon \(2021\)](#).

With the above definitions, we are now ready to introduce our estimator $\widehat{\mathbf{S}}$ of the noise covariance \mathbf{S} .

Definition 2.2 (Noise covariance estimator). Let $\mathbf{F} = \mathbf{Y} - \mathbf{X}\widehat{\mathbf{B}}$ and $\widehat{\mathbf{A}}$ be defined as before, we define our new noise covariance estimator as

$$\widehat{\mathbf{S}} = (n\mathbf{I}_T - \widehat{\mathbf{A}})^{-1} \left[\mathbf{F}^\top ((p+n)\mathbf{I}_n - \mathbf{X}\Sigma^{-1}\mathbf{X}^\top) \mathbf{F} - \widehat{\mathbf{A}}\mathbf{F}^\top \mathbf{F} - \mathbf{F}^\top \mathbf{F} \widehat{\mathbf{A}} \right] (n\mathbf{I}_T - \widehat{\mathbf{A}})^{-1}. \quad (11)$$

The construction of above $\widehat{\mathbf{S}}$ follows from novel applications of variants of Stein’s formulae in order to remove the bias in the naive estimator $\frac{1}{n}\mathbf{F}^\top \mathbf{F}$. The detailed derivation is presented in the Proof of Theorem 3.3. in Appendix B.2. The estimator $\widehat{\mathbf{S}}$ generalizes the scalar estimator (5) to the multi-task setting in the sense that for $T = 1$, $\widehat{\mathbf{S}}$ is exactly equal to (5). Note that unlike in (5), here $\mathbf{F}^\top \mathbf{F}$, $\widehat{\mathbf{A}}$ and $(n\mathbf{I}_T - \widehat{\mathbf{A}})$ are matrices of size $T \times T$: the order of matrix multiplication in $\widehat{\mathbf{S}}$ matters and should not be switched. This non-commutativity is not present for $T = 1$ in (5) where matrices in $\mathbb{R}^{T \times T}$ are reduced to scalars. Another special case of $\widehat{\mathbf{S}}$ can be seen in [Celentano and Montanari \(2021\)](#) for $T = 2$ where the matrix $\widehat{\mathbf{A}} \in \mathbb{R}^{2 \times 2}$ is diagonal and the two columns of $\mathbf{B} \in \mathbb{R}^{p \times 2}$ are two Lasso or Ridge estimators computed independently of each other, one for each task. Except in these two special cases — (5) for $T = 1$, [Celentano and Montanari \(2021\)](#) for $T = 2$ and two Lasso/Ridge — we are not aware of previously proposed estimators of the same form as $\widehat{\mathbf{S}}$.

The definition of our estimator $\widehat{\mathbf{S}}$ involves simple algebraic operations between the matrices \mathbf{X} , Σ , the residual \mathbf{F} and the interaction matrix $\widehat{\mathbf{A}}$. The multi-task estimate $\widehat{\mathbf{B}}$ in (6) can be efficiently solved by existing solver (e.g., `sklearn.linear_model.MultiTaskElasticNet` in scikit-learn library ([Pedregosa et al., 2011](#))), computation of \mathbf{F} is then straightforward, and computing the matrix $\widehat{\mathbf{A}}$ only requires inverting a matrix of size $|\mathcal{J}|$ thanks to the Sherman-Morrison-Woodbury formula ([Bellec and Romon, 2021](#), Section 5).

3. Theoretical analysis

3.1. Oracle and method-of-moments estimator

Before moving on to the theoretical analysis of $\widehat{\mathbf{S}}$, we state our randomness assumptions for \mathbf{E} , \mathbf{X} and we study two preliminary estimators: the oracle $\frac{1}{n}\mathbf{E}^\top \mathbf{E}$ and another estimator obtained by the method of moments.

Assumption 1 (Gaussian noise). $\mathbf{E} \in \mathbb{R}^{n \times T}$ is a Gaussian noise matrix with i.i.d. $\mathcal{N}_T(\mathbf{0}, \mathbf{S})$ rows, where $\mathbf{S} \in \mathbb{R}^{T \times T}$ is an unknown positive semi-definite matrix.

An oracle with access to the noise matrix \mathbf{E} may compute the oracle estimator $\widehat{\mathbf{S}}_{(\text{oracle})} \stackrel{\text{def}}{=} \frac{1}{n}\mathbf{E}^\top \mathbf{E}$, with convergence rate given by the following theorem, which will serve as a benchmark.

Proposition 3.1 (Convergence rate of $\widehat{\mathbf{S}}_{(\text{oracle})}$). *Under Assumption 1,*

$$\mathbb{E}[\|\widehat{\mathbf{S}}_{(\text{oracle})} - \mathbf{S}\|_{\text{F}}^2] = \frac{1}{n}[(\text{Tr}(\mathbf{S}))^2 + \text{Tr}(\mathbf{S}^2)]. \quad (12)$$

Consequently, $n^{-1}(\text{Tr}(\mathbf{S}))^2 \leq \mathbb{E}[\|\widehat{\mathbf{S}}_{(\text{oracle})} - \mathbf{S}\|_{\text{F}}^2] \leq 2n^{-1}(\text{Tr}(\mathbf{S}))^2$.

The next assumption concerns the design matrix \mathbf{X} with rows $\mathbf{x}_1^\top, \dots, \mathbf{x}_n^\top$.

Assumption 2 (Gaussian design). $\mathbf{X} \in \mathbb{R}^{n \times p}$ is a Gaussian design matrix with i.i.d. $\mathcal{N}_p(\mathbf{0}, \mathbf{\Sigma})$ rows, where $\mathbf{\Sigma}$ is a known positive definite matrix. The matrices \mathbf{E} and \mathbf{X} are independent.

Under the preceding assumptions, we obtain the following method-of-moments estimator, which extends the estimator for noise variance in [Dicker \(2014\)](#) to the multi-task setting. Its error will also serve as a benchmark.

Proposition 3.2. *Under Assumptions 1 and 2, the method-of-moments estimator defined as*

$$\widehat{\mathbf{S}}_{(\text{mm})} = \frac{(n+1+p)}{n(n+1)} \mathbf{Y}^\top \mathbf{Y} - \frac{1}{n(n+1)} \mathbf{Y}^\top \mathbf{X} \mathbf{\Sigma}^{-1} \mathbf{X}^\top \mathbf{Y} \quad (13)$$

is unbiased for \mathbf{S} , i.e., $\mathbb{E}[\widehat{\mathbf{S}}_{(\text{mm})}] = \mathbf{S}$. Furthermore, the Frobenius error is bounded from below as

$$\mathbb{E}[\|\widehat{\mathbf{S}}_{(\text{mm})} - \mathbf{S}\|_{\text{F}}^2] \geq \frac{p-2}{(n+1)^2} [\text{Tr}(\mathbf{S}) + \|\mathbf{\Sigma}^{\frac{1}{2}} \mathbf{B}^*\|_{\text{F}}^2]^2. \quad (14)$$

By (14), a larger norm $\|\mathbf{\Sigma}^{1/2} \mathbf{B}^*\|_{\text{F}}$ induces a larger variance for $\widehat{\mathbf{S}}_{(\text{mm})}$. Our goal with an estimate $\widehat{\mathbf{S}}$, when a good estimator $\widehat{\mathbf{B}}$ of the nuisance is available, is to improve upon the right-hand side of (14) when the estimation error $\|\mathbf{\Sigma}^{1/2}(\widehat{\mathbf{B}} - \mathbf{B}^*)\|_{\text{F}}$ is smaller than $\|\mathbf{\Sigma}^{1/2} \mathbf{B}^*\|_{\text{F}}$.

A high-probability upper bound of the form $\|\widehat{\mathbf{S}}_{(\text{mm})} - \mathbf{S}\|_{\text{F}}^2 \leq C \frac{n+p}{n^2} [\text{Tr}(\mathbf{S}) + \|\mathbf{\Sigma}^{\frac{1}{2}} \mathbf{B}^*\|_{\text{F}}^2]^2$, that matches the lower bound (14) when $p > n$, is a consequence of our main result Theorem 3.4 in the next subsection. Indeed, when $\widehat{\mathbf{B}} = \mathbf{0}$ then $\widehat{\mathbf{A}} = \mathbf{0}$ and our estimator $\widehat{\mathbf{S}}$ from Definition 2.2 coincides with $\widehat{\mathbf{S}}_{(\text{mm})}$ up to the minor modification of replacing $n+1$ by n in (13). This replacement is immaterial compared to the right-hand side in (14). Furthermore, such $\widehat{\mathbf{S}}$ corresponds to one of τ or λ being $+\infty$ in (6) and the aforementioned upper bound follows by taking $\tau = +\infty$ in the proof of Theorem 3.4 below. The empirical results in Section 4 confirm that $\widehat{\mathbf{S}}$ has smaller variance compared to $\widehat{\mathbf{S}}_{(\text{mm})}$ in simulations.

3.2. Theoretical results for proposed estimator $\widehat{\mathbf{S}}$

We have established lower bounds for the oracle estimator and the method-of-moments estimator that will serve as benchmarks. We turn to the analysis of the proposed estimator $\widehat{\mathbf{S}}$ from Definition 2.2 under the following additional assumptions.

Assumption 3 (High-dimensional regime). n, p satisfy $p/n \leq \gamma$ for a constant $\gamma \in (0, \infty)$.

For asymptotic statements such as those involving the stochastically bounded notation $O_p(\cdot)$ or the convergence in probability in (23) below, we implicitly consider a sequence of multi-task problems

indexed by n where $p, T, \mathbf{B}^*, \widehat{\mathbf{B}}, \mathbf{S}$ all implicitly depend on n . The assumptions, such as $p/n \leq \gamma$ above, are required to hold at all points of the sequence. In particular, $p/n \rightarrow \gamma'$ is allowed for any limit $\gamma' \leq \gamma$ under Assumption 3, although our results do not require a specific value for the limit.

Assumption 4. Recall τ is a regularization parameter in problem (6), we assume either one of the following:

- i) $\tau > 0$, and let $\tau' = \tau / \|\Sigma\|_{\text{op}}$.
- ii) $\tau = 0$ and for $c > 0$, $\mathbb{P}(U_1) \geq 1 - \frac{1}{T}$ and $\mathbb{P}(U_1) \rightarrow 1$ as $n \rightarrow \infty$, where $U_1 = \{\|\widehat{\mathbf{B}}\|_0 \leq n(1-c)/2\}$ is the event that $\widehat{\mathbf{B}}$ has at most $n(1-c)/2$ nonzero rows. Finally, $T \leq e^{\sqrt{n}}$.
- iii) $\tau = 0$, $\Sigma_{jj} = 1$ for all $j \in [p]$, $0 < \phi_* \leq \phi_{\min}(\Sigma) \leq \phi_{\max}(\Sigma) \leq \phi^*$, $T \leq e^{\sqrt{n}}$, $\|\mathbf{B}^*\|_0 \leq \min\{c^*(\gamma, \Sigma), c^{**}(\gamma, \Sigma)\}n$, $\lambda \geq \mu^*(\gamma, \Sigma)\sqrt{\text{Tr}(\mathbf{S})/n}$, where $\mu^*(\gamma, \Sigma) = (30\|\Sigma\|_{\text{op}})^{1/2}(2 + \sqrt{\gamma})$, $c^*(\gamma) = \sup_{c \in [0, \frac{1}{16} \wedge \gamma]} \{c \log(e\gamma/c) \leq 1/64\}$, $c^*(\gamma, \Sigma) = \frac{c^*(\gamma)\phi_{\min}(\Sigma)}{64}$, $c^{**}(\gamma, \Sigma) = \frac{\phi_{\min}(\Sigma)}{192\|\Sigma\|_{\text{op}}(2+\sqrt{\gamma})^2}$.

Assumption 4(i) requires that the Ridge penalty in (6) be enforced, so that the objective function is strongly convex. Assumption 4(ii), on the other hand, does not require strong convexity but that the number of nonzero rows of $\widehat{\mathbf{B}}$ is small enough with high-probability, which is a reasonable assumption when the tuning parameter λ in (6) is large enough and \mathbf{B}^* is sparse enough. While we do not prove in the present paper that $\mathbb{P}(U_1) \rightarrow 1$ under assumptions on the tuning parameter λ and the sparsity of \mathbf{B}^* , results of a similar nature have been obtained previously in several group-Lasso settings (Lounici et al., 2011, Theorem 3.1), (Liu and Zhang, 2009, Lemma 6), (Bellec and Romon, 2021, Lemma C.3), (Bellec and Kuchibhotla, 2019, Proposition 3.7). However, the proofs in those papers are not valid for the proportional regime. Under Assumption 4(iii), the following Proposition 3.3 provides a bound for the support size of $\widehat{\mathbf{B}}$ and shows that Assumption 4(iii) implies Assumption 4(ii).

Proposition 3.3. Under Assumption 4(iii), we have $\mathbb{P}(\|\widehat{\mathbf{B}}\|_0 \leq n/3) \geq 1 - \exp(-c(\gamma, \Sigma)n)$ for some positive constant $c(\gamma, \Sigma)$ depending on γ, Σ only.

Theorem 3.4. Suppose that Assumptions 1 to 4 hold for all n, p as $n \rightarrow \infty$, then almost surely

$$\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)(\widehat{\mathbf{S}} - \mathbf{S})(\mathbf{I}_T - \widehat{\mathbf{A}}/n)\|_{\text{F}} \leq \Theta_1 n^{-\frac{1}{2}} (\|\mathbf{F}\|_{\text{F}}^2/n + \|\mathbf{H}\|_{\text{F}}^2 + \text{Tr}(\mathbf{S})) \quad (15)$$

for some non-negative random variable Θ_1 of constant order, in the sense that $\mathbb{E}[\Theta_1^2] \leq C(\tau')(T \wedge (1 + \frac{p}{n}))(1 + \frac{p}{n}) \leq C(\gamma, \tau')$ under Assumption 4(i), and $\mathbb{E}[I(\Omega)\Theta_1^2] \leq C(\gamma, c)$ under Assumption 4(ii), where $I(\Omega)$ is the indicator function of an event Ω with $\mathbb{P}(\Omega) \rightarrow 1$.

Above, $\Theta_1 \geq 0$ is said to be of constant order because $\Theta_1 = O_P(1)$ follows from $\mathbb{E}[\Theta_1^2] \leq C(\gamma, \tau')$ or from $\mathbb{E}[I(\Omega)\Theta_1^2] \leq C(\gamma, c)$ if the stochastically bounded notation $O_P(1)$ is allowed to hide constants depending on (γ, τ') or (γ, c) only. In the left-hand side of (15), multiplication by $\mathbf{I}_T - \widehat{\mathbf{A}}/n$ on both sides of the error $\widehat{\mathbf{S}} - \mathbf{S}$ can be further removed, as

$$\|\widehat{\mathbf{S}} - \mathbf{S}\|_{\text{F}} \leq \|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)(\widehat{\mathbf{S}} - \mathbf{S})(\mathbf{I}_T - \widehat{\mathbf{A}}/n)\|_{\text{F}} \|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1}\|_{\text{op}}^2 \quad (16)$$

and the fact that $\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1}\|_{\text{op}}$ is bounded from above with high probability by a constant depending on γ, τ', c only. Upper bounds on $\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1}\|_{\text{op}}$ are formally stated in Appendix B.1.

We are now ready to present our main result on the error bounds for $\widehat{\mathbf{S}}$.

Theorem 3.5. *Let Assumptions 1 to 4 be fulfilled and $T = o(n)$. Then*

$$\|\widehat{\mathbf{S}} - \mathbf{S}\|_{\text{F}} \leq O_P(n^{-\frac{1}{2}})(\|\mathbf{F}\|_{\text{F}}^2/n), \quad (17)$$

$$\|\widehat{\mathbf{S}} - \mathbf{S}\|_{\text{F}} \leq O_P(n^{-\frac{1}{2}})[\text{Tr}(\mathbf{S}) + \|\mathbf{H}\|_{\text{F}}^2]. \quad (18)$$

Here the $O_P(n^{-\frac{1}{2}})$ notation involves constants depending on γ, τ', c .

It is instructive at this point to compare (18) with the lower bound (14) on the Frobenius error of the method-of-moments estimator. When $p \geq n$ then $\mathbb{E}[\|\widehat{\mathbf{S}}_{(\text{mm})} - \mathbf{S}\|_{\text{F}}^2] \geq \frac{c}{n}[\text{Tr}[\mathbf{S}] + \|\boldsymbol{\Sigma}^{1/2}\mathbf{B}^*\|_{\text{F}}^2]^2$; this is the situation where the Statistician does not attempt to estimate \mathbf{B}^* , and pays a price of $[\text{Tr}[\mathbf{S}] + \|\boldsymbol{\Sigma}^{1/2}\mathbf{B}^*\|_{\text{F}}^2]^2/n$. On the other hand, by definition of \mathbf{H} in (9), the right-hand side of (18), when squared, is of order $n^{-1}[\text{Tr}[\mathbf{S}] + \|\boldsymbol{\Sigma}^{1/2}(\widehat{\mathbf{B}} - \mathbf{B}^*)\|_{\text{F}}^2]^2$. Here the error bound only depends on \mathbf{B}^* through the estimation error for the nuisance $\|\boldsymbol{\Sigma}^{1/2}(\widehat{\mathbf{B}} - \mathbf{B}^*)\|_{\text{F}}^2$. This explains that when $\widehat{\mathbf{B}}$ is a good estimator of \mathbf{B}^* and $\|\boldsymbol{\Sigma}^{1/2}(\widehat{\mathbf{B}} - \mathbf{B}^*)\|_{\text{F}}^2$ is smaller compared to $\|\boldsymbol{\Sigma}^{1/2}\mathbf{B}^*\|_{\text{F}}^2$, the estimator $\widehat{\mathbf{S}}$ that leverages $\widehat{\mathbf{B}}$ will outperform the method-of-moments estimator $\widehat{\mathbf{S}}_{(\text{mm})}$ which does not attempt to estimate the nuisance \mathbf{B}^* .

Finally, the next results show that under additional assumptions, the estimator $\widehat{\mathbf{S}}$ enjoys Frobenius error bounds similar to the oracle estimator $\frac{1}{n}\mathbf{E}^\top \mathbf{E}$.

Assumption 5. $\text{SNR} \leq \text{snr}$ for some positive constant snr independent of n, p, T , where $\text{SNR} = \|\boldsymbol{\Sigma}^{\frac{1}{2}}\mathbf{B}^*\|_{\text{F}}^2/\text{Tr}(\mathbf{S})$ denotes the signal-to-noise ratio of the multi-task linear model (1).

Corollary 3.6. *Suppose that Assumptions 1, 2, 3, 4(i), 5 and $T = o(n)$ hold, then*

$$\|\widehat{\mathbf{S}} - \mathbf{S}\|_{\text{F}} \leq O_P(n^{-\frac{1}{2}}) \text{Tr}(\mathbf{S}), \quad (19)$$

where $O_P(\cdot)$ hides constants depending on $\gamma, \tau', \text{snr}$. Furthermore,

$$\begin{aligned} \|\widehat{\mathbf{S}} - \mathbf{S}\|_{\text{F}}^2 &\leq O_P(T/n) \|\mathbf{S}\|_{\text{F}}^2 = o_P(1) \|\mathbf{S}\|_{\text{F}}^2, \\ \|\|\widehat{\mathbf{S}}\|_* - \text{Tr}(\mathbf{S})\| &\leq O_P(\sqrt{T/n}) \text{Tr}(\mathbf{S}) = o_P(1) \text{Tr}(\mathbf{S}). \end{aligned}$$

Corollary 3.7. *Suppose that Assumptions 1, 2, 3, 4(ii) and $T = o(n)$ hold. If $\|\mathbf{B}^*\|_0 \leq (1-c)n/2$ and the tuning parameter λ is of the form $\lambda = \mu\sqrt{\text{Tr}(\mathbf{S})}/n$ for some positive constant μ , then*

$$\|\widehat{\mathbf{S}} - \mathbf{S}\|_{\text{F}} \leq O_P(n^{-\frac{1}{2}})(1 + \mu^2) \text{Tr}(\mathbf{S}), \quad (20)$$

where $O_P(\cdot)$ hides constants depending on $c, \gamma, \phi_{\min}(\boldsymbol{\Sigma})$.

The following corollary is a direct consequence of Corollary 3.7 and Proposition 3.3.

Corollary 3.8. *Suppose that Assumptions 1, 2, 3, 4(iii) and $T = o(n)$ hold. We have*

$$\|\widehat{\mathbf{S}} - \mathbf{S}\|_{\text{F}} \leq O_P(n^{-\frac{1}{2}}) \text{Tr}(\mathbf{S}), \quad (21)$$

where $O_P(\cdot)$ hides constants depending on $c, \gamma, \phi_{\min}(\boldsymbol{\Sigma}), \phi_{\max}(\boldsymbol{\Sigma})$.

Comparing Corollaries 3.6 to 3.8 with Proposition 3.1, we conclude that $\|\widehat{\mathbf{S}} - \mathbf{S}\|_{\text{F}}^2$ is of the same order as the Frobenius error of the oracle estimator in (12) up to constants depending on the signal-to-noise ratio, γ , and τ' under Assumption 4(i), and up to constants depending on $\mu, c, \gamma, \phi_{\min}(\mathbf{\Sigma})$ under Assumption 4(ii).

The error bounds in (18)–(20) are measured in Frobenius norm, similarly to existing works on noise covariance estimation Molstad (2022). Outside the context of linear regression models, much work has been devoted to covariance estimation in the operator norm. By the loose bound $\|\mathbf{M}\|_{\text{op}} \leq \|\mathbf{M}\|_{\text{F}}$, our upper bounds carry over to the operator norm. The same cannot be said for lower bounds, since for instance $\mathbb{E}[\|\widehat{\mathbf{S}}_{(\text{oracle})} - \mathbf{S}\|_{\text{op}}^2] \asymp n^{-1} \|\mathbf{S}\|_{\text{op}} \text{Tr}(\mathbf{S})$ (see, e.g., (Koltchinskii and Lounici, 2017, Corollary 2)).

3.3. Generalization error estimation

By analogy with single task models, we define the *generalization error* in multi-task models as the matrix $\mathbf{H}^\top \mathbf{H} + \mathbf{S}$ of size $T \times T$, whose (t, t') -th entry is $\mathbb{E}[(Y_t^{\text{new}} - \mathbf{x}_{\text{new}}^\top \widehat{\mathbf{B}} \mathbf{e}_t)(Y_{t'}^{\text{new}} - \mathbf{x}_{\text{new}}^\top \widehat{\mathbf{B}} \mathbf{e}_{t'}) | (\mathbf{X}, \mathbf{Y})]$ where $(Y_t^{\text{new}}, Y_{t'}^{\text{new}}, \mathbf{x}_{\text{new}})$ is independent of (\mathbf{X}, \mathbf{Y}) and has the same distribution as $(Y_{it}, Y_{it'}, \mathbf{x}_i)$ for some $i = 1, \dots, n$. Estimating the generalization error is useful for parameter tuning: since

$$\text{Tr}[\mathbf{H}^\top \mathbf{H} + \mathbf{S}] = \|\mathbf{\Sigma}^{1/2}(\widehat{\mathbf{B}} - \mathbf{B}^*)\|_{\text{F}}^2 + \text{Tr}[\mathbf{S}], \quad (22)$$

minimizing an estimator of $\text{Tr}[\mathbf{H}^\top \mathbf{H} + \mathbf{S}]$ is a useful proxy to minimize the Frobenius error $\|\mathbf{\Sigma}^{1/2}(\widehat{\mathbf{B}} - \mathbf{B}^*)\|_{\text{F}}^2$ of $\widehat{\mathbf{B}}$.

The following theorem suggests an estimate for the generalization error matrix as well as a consistent estimator for its trace (22).

Theorem 3.9 (Generalization error). *Let Assumptions 1 to 4 be fulfilled. Then*

$$\|\mathbf{F}^\top \mathbf{F}/n - (\mathbf{I}_T - \widehat{\mathbf{A}}/n)(\mathbf{H}^\top \mathbf{H} + \mathbf{S})(\mathbf{I}_T - \widehat{\mathbf{A}}/n)\|_{\text{F}} \leq \Theta_2 n^{-\frac{1}{2}} (\|\mathbf{F}\|_{\text{F}}^2/n + \|\mathbf{H}\|_{\text{F}}^2 + \text{Tr}(\mathbf{S})),$$

for some non-negative random variable Θ_2 of constant order, in the sense that $\mathbb{E}[\Theta_2] \leq C(\gamma, \tau')$ under Assumption 4(i), and with $\mathbb{E}[I(\Omega)\Theta_2] \leq C(\gamma, c)$ under Assumption 4(ii) or (iii), where $I(\Omega)$ is the indicator function of an event Ω with $\mathbb{P}(\Omega) \rightarrow 1$.

Furthermore, if $T = o(n)$ as $n, p \rightarrow \infty$ while τ', γ, c stay constant, then

$$\frac{\text{Tr}(\mathbf{S}) + \|\mathbf{H}\|_{\text{F}}^2}{\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1} \mathbf{F}^\top\|_{\text{F}}^2/n} \xrightarrow{p} 1. \quad (23)$$

In the above theorem, \mathbf{S} and \mathbf{H} are unknown, while $\widehat{\mathbf{A}}$ and \mathbf{F} can be computed from the observed data (\mathbf{X}, \mathbf{Y}) . Thus (23) shows that $\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1} \mathbf{F}^\top\|_{\text{F}}^2/n$ is a consistent estimate for the unobservable quantity $\text{Tr}(\mathbf{S}) + \|\mathbf{H}\|_{\text{F}}^2$ without requiring the knowledge of $\mathbf{\Sigma}$. Theorem 3.9 also suggests an useful method for selecting tuning parameter, i.e., to choose the parameter that minimizes the estimated generalization error.

3.4. Out-of-sample error estimation

In this section, we present a by-product of our techniques for estimating the noise covariance. For evaluating the performance of a regression method on a new data, we define the out-of-sample error for

the multi-task linear model (1) as

$$\mathbb{E}[(\widehat{\mathbf{B}} - \mathbf{B}^*)^\top \mathbf{x}_{\text{new}} \mathbf{x}_{\text{new}}^\top (\widehat{\mathbf{B}} - \mathbf{B}^*) | (X, Y)] = \mathbf{H}^\top \mathbf{H},$$

where \mathbf{x}_{new} is independent of the data $(X; Y)$ with the same distribution as any row of X . The following theorem on estimation of out-of-sample error is a by-product of our technique for constructing $\widehat{\mathbf{S}}$.

Theorem 3.10 (Out-of-sample error). *Under the same conditions of Theorem 3.4, with $\mathbf{Z} = \mathbf{X}\Sigma^{-\frac{1}{2}}$, we have*

$$\begin{aligned} & \|(\mathbf{I}_T - \widehat{\mathbf{A}}/n) \mathbf{H}^\top \mathbf{H} (\mathbf{I}_T - \widehat{\mathbf{A}}/n) - \frac{1}{n^2} (\mathbf{F}^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{F} + \widehat{\mathbf{A}} \mathbf{F}^\top \mathbf{F} + \mathbf{F}^\top \mathbf{F} \widehat{\mathbf{A}} - p \mathbf{F}^\top \mathbf{F})\|_{\text{F}} \\ & \leq \Theta_3 n^{-\frac{1}{2}} (\|\mathbf{H}\|_{\text{F}}^2 + \|\mathbf{F}\|_{\text{F}}^2/n) \end{aligned}$$

for some non-negative random variable Θ_3 of constant order, in the sense that $\mathbb{E}[\Theta_3] \leq C(\gamma, \tau')$ under Assumption 4(i), and with $\mathbb{E}[I(\Omega)\Theta_3] \leq C(\gamma, c)$ under Assumption 4(ii) or (iii), where $I(\Omega)$ is the indicator function of an event Ω with $\mathbb{P}(\Omega) \rightarrow 1$.

Theorem 3.10 generalizes the result in Bellec (2020) to multi-task setting. While the out-of-sample error $\mathbf{H}^\top \mathbf{H}$ is unknown, the quantities \mathbf{Z} , \mathbf{F} , $\widehat{\mathbf{A}}$ are observable when Σ is known. Since typically the quantity $(\|\mathbf{H}\|_{\text{F}}^2 + \|\mathbf{F}\|_{\text{F}}^2/n)$ is of a constant order, Theorem 3.10 suggests the following estimate of $\mathbf{H}^\top \mathbf{H}$:

$$\frac{1}{n^2} (\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1} (\mathbf{F}^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{F} + \widehat{\mathbf{A}} \mathbf{F}^\top \mathbf{F} + \mathbf{F}^\top \mathbf{F} \widehat{\mathbf{A}} - p \mathbf{F}^\top \mathbf{F}) (\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1},$$

which can further be used for parameter tuning in multi-task linear models.

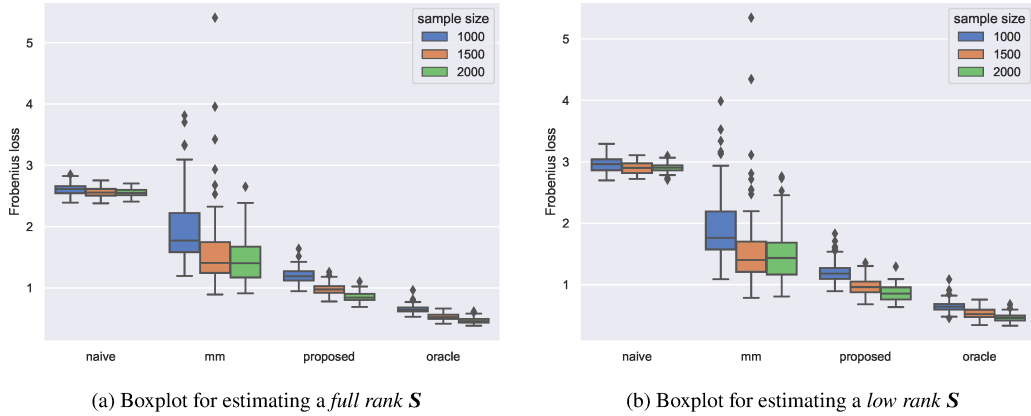


Figure 1: Boxplots for Frobenius norm loss over 100 repetitions.

Table 1. Frobenius norm loss for different methods

		$n = 1000$		$n = 1500$		$n = 2000$	
S	method	mean	sd	mean	sd	mean	sd
full rank	naive	2.610	0.086	2.560	0.075	2.555	0.062
	mm	1.970	0.549	1.586	0.621	1.475	0.381
	proposed	1.206	0.114	0.982	0.098	0.853	0.081
	oracle	0.652	0.061	0.534	0.052	0.469	0.045
low rank	naive	2.960	0.121	2.899	0.100	2.901	0.077
	mm	1.916	0.559	1.560	0.643	1.480	0.428
	proposed	1.208	0.178	0.971	0.130	0.861	0.120
	oracle	0.654	0.096	0.531	0.081	0.464	0.065

4. Numerical experiments

In this section, we evaluate the empirical performance of our proposed method and compare it with relevant methods for noise covariance estimation.

Regarding parameters for our simulations, we set $T = 20$, $p = 1.5n$ and n equals successively 1000, 1500, 2000. We consider two types of noise covariance matrix: (i) S is full-rank with (t, t') -th entry $S_{t,t'} = \frac{\cos(t-t')}{1+\sqrt{|t-t'|}}$; (ii) S is low-rank with $S = \mathbf{u}\mathbf{u}^\top$, where $\mathbf{u} \in \mathbb{R}^{T \times 10}$ has i.i.d. entries from $\mathcal{N}(0, 1/T)$. The design matrix $X \in \mathbb{R}^{n \times p}$ is constructed by independently sampling its rows from $\mathcal{N}_p(\mathbf{0}, \Sigma)$ with $\Sigma_{jk} = 0.5^{|j-k|}$. To build the coefficient matrix \mathbf{B}^* , we first set its sparsity pattern, i.e., we define the support \mathcal{S} with cardinality $|\mathcal{S}| = 0.1p$, then we generate an intermediate matrix $\mathbf{B} \in \mathbb{R}^{p \times T}$. The j -th row of \mathbf{B} is sampled from $\mathcal{N}_T(\mathbf{0}, p^{-1}\mathbf{I}_T)$ if $j \in \mathcal{S}$, otherwise we set it to be the zero vector. Finally we let $\mathbf{B}^* = \mathbf{B}[\text{Tr}(\mathbf{S})/\text{Tr}(\mathbf{B}^\top \Sigma \mathbf{B})]^\frac{1}{2}$, which forces a signal-to-noise ratio of exactly 1.

For calculation of multi-task estimates $\hat{\mathbf{B}}$ in (6), we use Python library Scikit-learn (Pedregosa et al., 2011). More precisely we invoke `MultiTaskElasticNetCV` to obtain $\hat{\mathbf{B}}$ by 5-fold cross-validation with parameters `l1-ratio=[0.5, 0.7, 0.9, 1]`, `n_alpha=100`. To compute the interaction matrix $\hat{\mathbf{A}}$ we use the efficient implementation described in (Bellec and Romon, 2021, Section 5). The full code needed to reproduce our experiments is provided in the supplementary material Tan, Romon and Bellec (2023).

We compare our proposed estimator $\hat{\mathbf{S}}$ (11) with relevant estimators including (1) the naive estimate $\hat{\mathbf{S}}_{(\text{naive})} = n^{-1}\mathbf{F}^\top \mathbf{F}$, (2) the method-of-moments estimate $\hat{\mathbf{S}}_{(\text{mm})}$ defined in Proposition 3.2, and (3) the oracle estimate $\hat{\mathbf{S}}_{(\text{oracle})} = n^{-1}\mathbf{E}^\top \mathbf{E}$. The performance of each estimator is measured in Frobenius norm error: for instance, $\|\hat{\mathbf{S}} - \mathbf{S}\|_F$ is the error for proposed estimator $\hat{\mathbf{S}}$. For each configurations of (n, p) and the aforementioned two types of noise covariance S , we run 100 repetitions and report in Figure 1 the boxplots of the Frobenius error from different methods, and report in Table 1 the corresponding mean and standard deviation of the Frobenius error.

Figure 1 and Table 1 show that, besides the oracle estimator, our proposed estimator has the best performance with significantly smaller loss compared to the naive and method-of-moments estimators. As expected, oracle estimator outperforms all three other methods. However, since it is typically not available in practice, we recommend using our proposed estimator because it has the nearest performance to oracle estimator.

Since the estimation target S is a $T \times T$ matrix, we also want to compare different estimators in terms of the bias and standard deviation for each entry of S . Figure 2 and Figure 3 present the heatmaps of bias and standard deviation from different estimators for full-rank and low-rank S when $n = 1000$. We display

the remaining heatmaps from other n (1500 and 2000) in the supplementary material [Tan, Romon and Bellec \(2023\)](#) as they exhibit similar pattern to Figures 2 and 3.

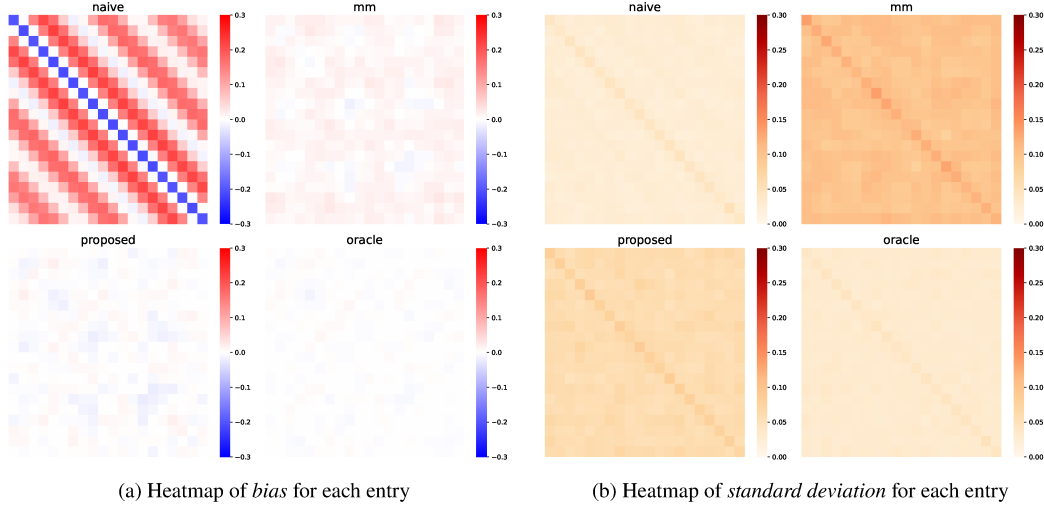


Figure 2: Heatmaps for estimation of full rank S with $n = 1000$ over 100 repetitions.

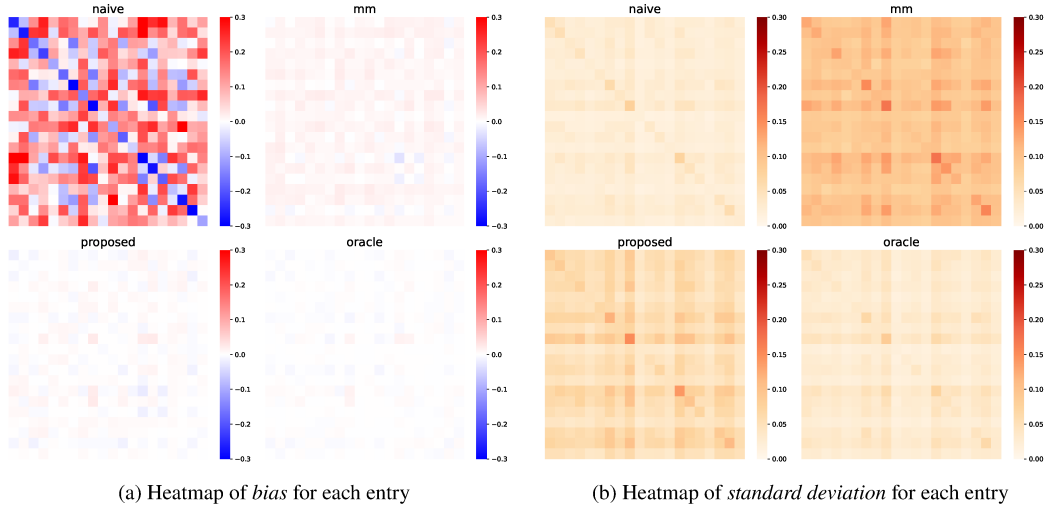


Figure 3: Heatmaps for estimation of low rank S with $n = 1000$ over 100 repetitions.

The comparison pattern of the four estimation methods in Figures 2 and 3 are consistent with Figure 1. It is not surprising that the oracle estimator has smallest bias and standard deviation. Figures 2 and 3 also

convince us that the naive estimator has large bias, though it has small standard deviation. The method-of-moments estimator is unbiased but its variance is relatively large, which means its performance is not stable, as was reflected in Figure 1 and Table 1. Our proposed estimator improves on both the naive and method-of-moments estimators because it has much smaller bias than the former, while having smaller standard deviation than the latter. Therefore, our proposed estimator is the clear winner besides the impractical oracle estimator.

5. Limitations and future work

One limitation of the proposed estimator \hat{S} is that its construction necessitates the knowledge of Σ . Let us first mention that the estimator $n^{-1}\|(I_T - \hat{A}/n)^{-1}F^\top\|_F$ of $\text{Tr}(S) + \|\Sigma^{1/2}(\hat{B} - B^*)\|_F^2$ in Theorem 3.9 does not require knowing Σ . Thus, this estimator can further be used as a proxy of the error $\|\Sigma^{1/2}(\hat{B} - B^*)\|_F^2$, say for parameter tuning, without the knowledge of Σ . The problem of estimating S with known Σ was studied in Celentano and Montanari (2021) for $T = 2$: in this inaccurate covariate model and for $p/n \leq \gamma$, our results yield the convergence rate $n^{-1/2}$ for S which improves upon the rate n^{-c_0} for a non-explicit constant $c_0 > 0$ in (Celentano and Montanari, 2021, Theorem 2.1).

In order to use \hat{S} when Σ is unknown, one may plug-in an estimator $\hat{\Sigma}$ in Equation (11), resulting in an extra term of order $\|\hat{\Sigma}^{-1} - \Sigma^{-1}\|_{\text{op}}\|F\|_F$ for the Frobenius error. See (Dicker, 2014, §4) for related discussions in the $T = 1$ (single-task) case. While, under the proportional regime $p/n \rightarrow \gamma$, no estimator is consistent for all covariance matrices Σ in operator norm, consistent estimators do exist under additional structural assumptions Bickel and Levina (2008), Cai, Zhang and Zhou (2010), El Karoui (2008). If available, additional unlabeled samples $(x_i)_{i \geq n+1}$ can also be used to construct norm-consistent estimator of Σ .

Future directions include extending estimator \hat{S} to utilize other estimators of the nuisance B^* than the multi-task elastic-net (6); for instance (7) or the estimators studied in Bertrand et al. (2019), Geer and Stucky (2016), Molstad (2022). In the simpler case where columns of B^* are estimated independently on each task, e.g., if the T columns of \hat{B} are Lasso estimators $(\hat{\beta}^{(t)})_{t \in [T]}$ each computed from $y^{(t)}$, then minor modifications of our proof yield that the estimator (11) with $\hat{A} = \text{diag}(\|\hat{\beta}^{(1)}\|_0, \dots, \|\hat{\beta}^{(T)}\|_0)$ enjoys similar Frobenius norm bounds of order $n^{-1/2}$.

Appendix A: Proof of Propositions 3.1 to 3.3

Notation. We first introduce a few notations that will be used throughout the proofs. We use indexes i and l only to loop or sum over $[n] = \{1, 2, \dots, n\}$, use j and k only to loop or sum over $[p] = \{1, 2, \dots, p\}$, use t and t' only to loop or sum over $[T] = \{1, 2, \dots, T\}$, so that e_i (and e_l) refer to the i -th (and l -th) canonical basis vector in \mathbb{R}^n , e_j (and e_k) refer to the j -th (and k -th) canonical basis vector in \mathbb{R}^p , e_t (and $e_{t'}$) refer to the t -th (and t' -th) canonical basis vector in \mathbb{R}^T . For any two real numbers a and b , let $a \vee b = \max(a, b)$, and $a \wedge b = \min(a, b)$. Positive constants that depend on γ, τ' only are denoted by $C(\gamma, \tau')$, and positive constants that depend on γ, c only are denoted by $C(\gamma, c)$. The values of these constants may vary from place to place.

Proof of Proposition 3.1. Let $S = \sum_{t=1}^T \sigma_t^2 u_t u_t^\top$ be the spectral decomposition of S , then $\|E^\top E - nS\|_F^2 = \sum_{t \in [T]} \sum_{t' \in [T]} [u_{t'}^\top (E^\top E - nS) u_t]^2$. We now compute the expectation of one term indexed by (t, t') . The random variable $u_{t'}^\top (E^\top E - nS) u_t$ is the sum of n i.i.d. mean zero random variables with the same distribution as $z_{t'} z_t - u_{t'}^\top S u_t$ where $(z_t, z_{t'}) \sim \mathcal{N}_2(\mathbf{0}, \text{diag}(\sigma_t^2, \sigma_{t'}^2))$. Thus

$$\mathbb{E}[(u_{t'}^\top (E^\top E - nS) u_t)^2] = n \text{Var}[z_{t'} z_t - u_{t'}^\top S u_t] = n(2\sigma_t^4 I_{t=t'} + n\sigma_t^2 \sigma_{t'}^2 I_{t \neq t'})$$

due to $\text{Var}[\chi_1^2] = 2$ if $t = t'$ and independence if $t \neq t'$. Summing over all $(t, t') \in [T] \times [T]$ yields $2n \sum_{t=1}^T \sigma_t^4 + n \sum_{t \neq t'} \sigma_t^2 \sigma_{t'}^2 = n \sum_{t=1}^T \sigma_t^4 + n (\sum_{t=1}^T \sigma_t^2)^2 = n \|\mathbf{S}\|_F^2 + n [\text{Tr}(\mathbf{S})]^2$ as desired.

The inequality simply follows from $\|\mathbf{S}\|_F^2 \leq [\text{Tr}(\mathbf{S})]^2$ since \mathbf{S} is positive semi-definite. ■

Proof of Proposition 3.2. Without loss of generality, we assume $\mathbf{\Sigma} = \mathbf{I}_p$. For general positive definite $\mathbf{\Sigma}$, the proof follows by replacing $(\mathbf{X}, \mathbf{B}^*)$ with $(\mathbf{X}\mathbf{\Sigma}^{-\frac{1}{2}}, \mathbf{\Sigma}^{\frac{1}{2}}\mathbf{B}^*)$.

We first derive the method-of-moments estimator $\widehat{\mathbf{S}}_{(\text{mm})}$. Under Assumptions 1 and 2 with $\mathbf{\Sigma} = \mathbf{I}_p$, \mathbf{X} has i.i.d. rows from $\mathcal{N}_p(\mathbf{0}, \mathbf{I}_p)$, \mathbf{E} has i.i.d. rows from $\mathcal{N}_T(\mathbf{0}, \mathbf{S})$, and \mathbf{X} and \mathbf{E} are independent. Then, the expectations of $\mathbf{Y}^\top \mathbf{Y}$ and $\mathbf{Y}^\top \mathbf{X} \mathbf{X}^\top \mathbf{Y}$ are given by

$$\mathbb{E}(\mathbf{Y}^\top \mathbf{Y}) = \mathbb{E}[(\mathbf{X}\mathbf{B}^* + \mathbf{E})^\top (\mathbf{X}\mathbf{B}^* + \mathbf{E})] = n(\mathbf{B}^{*\top} \mathbf{B}^* + \mathbf{S}), \quad (24)$$

and

$$\begin{aligned} \mathbb{E}(\mathbf{Y}^\top \mathbf{X} \mathbf{X}^\top \mathbf{Y}) &= \mathbb{E}[(\mathbf{X}\mathbf{B}^* + \mathbf{E})^\top \mathbf{X} \mathbf{X}^\top (\mathbf{X}\mathbf{B}^* + \mathbf{E})] \\ &= \mathbb{E}(\mathbf{B}^{*\top} \mathbf{X}^\top \mathbf{X} \mathbf{X}^\top \mathbf{X} \mathbf{B}^*) + \mathbb{E}(\mathbf{E}^\top \mathbf{X} \mathbf{X}^\top \mathbf{E}) \\ &= \mathbf{B}^{*\top} \mathbb{E}(\mathbf{X}^\top \mathbf{X} \mathbf{X}^\top \mathbf{X}) \mathbf{B}^* + \mathbb{E}(\mathbf{E}^\top \mathbf{X} \mathbf{X}^\top \mathbf{E}) \\ &= n(n + p + 1) \mathbf{B}^{*\top} \mathbf{B}^* + np\mathbf{S}, \end{aligned} \quad (25)$$

where the last line uses

$$\begin{aligned} &\mathbb{E}(\mathbf{X}^\top \mathbf{X} \mathbf{X}^\top \mathbf{X}) \\ &= \mathbb{E}\left[\sum_{i=1}^n (\mathbf{x}_i \mathbf{x}_i^\top) \sum_{l=1}^n (\mathbf{x}_l \mathbf{x}_l^\top)\right] \\ &= \sum_{i \neq l} \mathbb{E}(\mathbf{x}_i \mathbf{x}_i^\top \mathbf{x}_l \mathbf{x}_l^\top) + \sum_{i=l} \mathbb{E}(\mathbf{x}_i \mathbf{x}_i^\top \mathbf{x}_l \mathbf{x}_l^\top) \\ &= n(n-1) \mathbf{I}_p^2 + n \mathbb{E}(\mathbf{x}_1 \mathbf{x}_1^\top \mathbf{x}_1 \mathbf{x}_1^\top) \\ &= n(n-1) \mathbf{I}_p + n[2\mathbf{I}_p^2 + \text{Tr}(\mathbf{I}_p) \mathbf{I}_p] \\ &= n(n + p + 1) \mathbf{I}_p, \end{aligned}$$

and

$$\mathbb{E}(\mathbf{E}^\top \mathbf{X} \mathbf{X}^\top \mathbf{E}) = \mathbb{E}[\mathbb{E}(\mathbf{E}^\top \mathbf{X} \mathbf{X}^\top \mathbf{E} | \mathbf{E})] = \mathbb{E}[\mathbf{E}^\top \mathbb{E}(\mathbf{X} \mathbf{X}^\top) \mathbf{E}] = np\mathbf{S}.$$

Solving for \mathbf{S} from the system of equations (24) and (25), we obtain the method-of-moments estimator

$$\widehat{\mathbf{S}}_{(\text{mm})} = \frac{(n + p + 1)}{n(n + 1)} \mathbf{Y}^\top \mathbf{Y} - \frac{1}{n(n + 1)} \mathbf{Y}^\top \mathbf{X} \mathbf{X}^\top \mathbf{Y},$$

and $\mathbb{E}[\widehat{\mathbf{S}}_{(\text{mm})}] = \mathbf{S}$.

Now we derive the variance lower bound for $\widehat{\mathbf{S}}_{(\text{mm})}$. Since $\mathbb{E}[\widehat{\mathbf{S}}_{(\text{mm})}] = \mathbf{S}$, $\mathbb{E}[\|\widehat{\mathbf{S}}_{(\text{mm})} - \mathbf{S}\|_F^2] = \sum_{t, t'} \text{Var}\{\widehat{\mathbf{S}}_{(\text{mm})}[t, t']\}$. By definition of $\widehat{\mathbf{S}}_{(\text{mm})}$,

$$[\widehat{\mathbf{S}}_{(\text{mm})}]_{t, t'} = \frac{n + p + 1}{n(n + 1)} [\mathbf{y}^{(t)}]^\top \mathbf{y}^{(t')} - \frac{1}{n(n + 1)} [\mathbf{y}^{(t)}]^\top \mathbf{X} \mathbf{\Sigma}^{-1} \mathbf{X}^\top \mathbf{y}^{(t')}.$$

Since $\mathbf{y}^{(t)} = \mathbf{X}\boldsymbol{\beta}^{(t)} + \boldsymbol{\varepsilon}^{(t)}$, $\mathbf{y}^{(t')} = \mathbf{X}\boldsymbol{\beta}^{(t')} + \boldsymbol{\varepsilon}^{(t')}$, for $t \neq t'$, without loss of generality, we assume $\boldsymbol{\beta}^{(t)} = a_0\mathbf{e}_1$ and $\boldsymbol{\beta}^{(t')} = a_1\mathbf{e}_1 + a_2\mathbf{e}_2$ for some constants a_0, a_1, a_2 . If necessary, we could let $\mathbf{u}_1 = \boldsymbol{\beta}^{(t)} / \|\boldsymbol{\beta}^{(t)}\|$, and $\mathbf{u}_2 = \tilde{\mathbf{u}}_2 / \|\tilde{\mathbf{u}}_2\|$ where $\tilde{\mathbf{u}}_2 = \boldsymbol{\beta}^{(t')} - \mathbf{P}_{\mathbf{u}_1}\boldsymbol{\beta}^{(t')}$, and completing the basis to obtain an orthonormal basis $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_p\}$ for \mathbb{R}^p . Let $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_p]$, then \mathbf{U} is an orthogonal matrix, hence $\mathbf{X}\mathbf{U}$ and \mathbf{X} have the same distribution, only the first coordinate of $\mathbf{U}^\top \boldsymbol{\beta}^{(t)}$ is nonzero, and only the first two coordinates of $\mathbf{U}^\top \boldsymbol{\beta}^{(t')}$ are nonzero. That is, we could perform change of variables by replacing $(\mathbf{X}, \boldsymbol{\beta}^{(t)}, \boldsymbol{\beta}^{(t')})$ with $(\mathbf{X}\mathbf{U}, \mathbf{U}^\top \boldsymbol{\beta}^{(t)}, \mathbf{U}^\top \boldsymbol{\beta}^{(t')})$.

Therefore, $\mathbf{y}^{(t)}$ and $\mathbf{y}^{(t')}$ are independent of $\{\mathbf{X}\mathbf{e}_j : 3 \leq j \leq p\}$. Let $\mathcal{F} = \sigma(\mathbf{y}^{(t)}, \mathbf{y}^{(t')}, \mathbf{X}\mathbf{e}_1, \mathbf{X}\mathbf{e}_2)$ be the σ -field generated by $(\mathbf{y}^{(t)}, \mathbf{y}^{(t')}, \mathbf{X}\mathbf{e}_1, \mathbf{X}\mathbf{e}_2)$, then

$$\text{Var}\{[\widehat{\mathbf{S}}_{(\text{mm})}]_{t,t'}\} \geq \mathbb{E}[\text{Var}\{[\widehat{\mathbf{S}}_{(\text{mm})}]_{t,t'} | \mathcal{F}\}] = \frac{1}{n^2(n+1)^2} \mathbb{E}[\text{Var}\{[\mathbf{y}^{(t)}]^\top \mathbf{X} \mathbf{X}^\top \mathbf{y}^{(t')} | \mathcal{F}\}].$$

Note that in the above display,

$$[\mathbf{y}^{(t)}]^\top \mathbf{X} \mathbf{X}^\top \mathbf{y}^{(t')} = \sum_{j=1}^2 [\mathbf{y}^{(t)}]^\top \mathbf{X} \mathbf{e}_j \mathbf{e}_j^\top \mathbf{X}^\top \mathbf{y}^{(t')} + \sum_{j=3}^p [\mathbf{y}^{(t)}]^\top \mathbf{X} \mathbf{e}_j \mathbf{e}_j^\top \mathbf{X}^\top \mathbf{y}^{(t')},$$

where the first term is measurable with respect to \mathcal{F} , and the second term is a quadratic form

$$\sum_{j=3}^p [\mathbf{y}^{(t)}]^\top \mathbf{X} \mathbf{e}_j \mathbf{e}_j^\top \mathbf{X}^\top \mathbf{y}^{(t')} = \sum_{j=3}^p \mathbf{e}_j^\top \mathbf{X}^\top \mathbf{y}^{(t')} [\mathbf{y}^{(t)}]^\top \mathbf{X} \mathbf{e}_j = \boldsymbol{\xi}^\top \boldsymbol{\Lambda} \boldsymbol{\xi},$$

here $\boldsymbol{\xi} = [\mathbf{e}_3^\top \mathbf{X}^\top, \dots, \mathbf{e}_p^\top \mathbf{X}^\top]^\top \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n(p-2)})$, and $\boldsymbol{\Lambda} = \mathbf{I}_{p-2} \otimes \mathbf{y}^{(t')} [\mathbf{y}^{(t)}]^\top$. Thus, for $t \neq t'$,

$$\begin{aligned} \text{Var}\{[\widehat{\mathbf{S}}_{(\text{mm})}]_{t,t'}\} &\geq \frac{1}{n^2(n+1)^2} \mathbb{E}\{\text{Var}\{\boldsymbol{\xi}^\top \boldsymbol{\Lambda} \boldsymbol{\xi} | \mathcal{F}\}\} \\ &= \frac{1}{n^2(n+1)^2} \mathbb{E}\{\|\boldsymbol{\Lambda}\|_{\text{F}}^2 + \text{Tr}(\boldsymbol{\Lambda}^2)\} \\ &\geq \frac{1}{n^2(n+1)^2} \mathbb{E}[\|\boldsymbol{\Lambda}\|_{\text{F}}^2] \\ &= \frac{p-2}{n^2(n+1)^2} \mathbb{E}[\|\mathbf{y}^{(t)}\|^2 \|\mathbf{y}^{(t')}\|^2]. \end{aligned}$$

For $t = t'$, using a similar argument we obtain

$$\text{Var}\{[\widehat{\mathbf{S}}_{(\text{mm})}]_{t,t'}\} \geq \frac{p-1}{n^2(n+1)^2} \mathbb{E}[\|\mathbf{y}^{(t)}\|^2 \|\mathbf{y}^{(t')}\|^2].$$

Summing over all $(t, t') \in [T] \times [T]$ yields

$$\begin{aligned} \mathbb{E}[\|\widehat{\mathbf{S}}_{(\text{mm})} - \mathbf{S}\|_{\text{F}}^2] &\geq \frac{p-2}{n^2(n+1)^2} \sum_{t,t'} \mathbb{E}[\|\mathbf{y}^{(t)}\|^2 \|\mathbf{y}^{(t')}\|^2] \\ &= \frac{p-2}{n^2(n+1)^2} \mathbb{E}[\|\mathbf{Y}\|_{\text{F}}^4] \\ &\geq \frac{p-2}{n^2(n+1)^2} (\mathbb{E}[\|\mathbf{Y}\|_{\text{F}}^2])^2 \end{aligned}$$

$$= \frac{p-2}{(n+1)^2} [\text{Tr}(\mathbf{S}) + \|\mathbf{B}^*\|^2]^2.$$

■

Proof of Proposition 3.3. Let us first introduce the following two lemmas and defer their proofs to the end of this section.

Lemma A.1 (deterministic bound for $\|\widehat{\mathbf{B}}\|_0$). *For the multi-task Lasso estimator*

$$\widehat{\mathbf{B}} = \arg \min_{\mathbf{B} \in \mathbb{R}^{p \times T}} \left(\frac{1}{2n} \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_{\text{F}}^2 + \lambda \|\mathbf{B}\|_{2,1} \right),$$

we have

$$\|\widehat{\mathbf{B}}\|_0 \leq \|\mathbf{X}\|_{\text{op}}^2 \max \left\{ \frac{2\|\mathbf{E}\|_{\text{F}}^2}{n^2 \lambda^2}, \frac{4\|\mathbf{B}^*\|_0}{n \kappa^2} \right\},$$

$$\text{where } \kappa^2 = \inf_{\mathbf{B} \in \mathbb{R}^{p \times T} : \|\mathbf{B}\|_{2,1} \leq \|\mathbf{B}^*\|_{2,1}} \frac{\|\mathbf{X}(\mathbf{B} - \mathbf{B}^*)\|_{\text{F}}^2}{n \|\mathbf{B} - \mathbf{B}^*\|_{\text{F}}^2}.$$

Lemma A.1 is a short deterministic argument that provides an upper bound for $\|\widehat{\mathbf{B}}\|_0$ of the multi-task Lasso. It follows the same lines as that for the Lasso in (Bellec, 2020, Section 10). This argument is sufficient in the $n \asymp p$ regime of the present paper; in the regime $s \ll n \ll p$ regime, sharper bounds are available ((Lounici et al., 2011, Theorem 3.1), (Liu and Zhang, 2009, Lemma 6), (Bellec and Romon, 2021, Lemma C.3), (Bellec and Kuchibhotla, 2019, Proposition 3.7)).

Lemma A.2 is a generalization of (Lecué and Mendelson, 2017, Lemma 2.7) to the multi-task setting and will be useful to provide a lower bound for κ in Lemma A.1.

Lemma A.2. *Let $\mathbf{X} \in \mathbb{R}^{n \times p}$, $k \in [p]$, $\delta > 0$ be such that*

$$\inf_{\mathbf{b} \in \mathbb{R}^p : \|\mathbf{b}\|_0 \leq k} \frac{\|\mathbf{X}\mathbf{b}\|}{\sqrt{n}\|\mathbf{b}\|} \geq \delta, \quad (26)$$

then for any matrix $\mathbf{A} \in \mathbb{R}^{p \times T}$ we have

$$\left(1 - \frac{1}{k}\right) \frac{\|\mathbf{X}\mathbf{A}\|_{\text{F}}^2}{n} \geq \delta^2 \left(1 - \frac{1}{k}\right) \|\mathbf{A}\|_{\text{F}}^2 - \frac{(\sum_{j=1}^p \|\mathbf{A}^\top \mathbf{e}_j\|)^2}{k} \left[\sum_j \left(\mu_j \frac{\|\mathbf{X}\mathbf{e}_j\|^2}{n} - \delta^2 \right) \right],$$

$$\text{where } \mu_j = \frac{\|\mathbf{A}^\top \mathbf{e}_j\|}{\sum_j \|\mathbf{A}^\top \mathbf{e}_j\|}.$$

Our first step is to bound κ from below using Lemma A.2. To this end, we first verify condition (26) by finding suitable expressions for k and δ . Note that

$$\inf_{\mathbf{b} \in \mathbb{R}^p : \|\mathbf{b}\|_0 \leq k} \frac{\|\mathbf{X}\mathbf{b}\|}{\sqrt{n}\|\mathbf{b}\|} \geq \phi_{\min}^{1/2}(\Sigma) \inf_{\mathbf{b} \in \mathbb{R}^p : \|\mathbf{b}\|_0 \leq k} \frac{\|\mathbf{Z}\Sigma^{1/2}\mathbf{b}\|}{\sqrt{n}\|\Sigma^{1/2}\mathbf{b}\|} = \phi_{\min}^{1/2}(\Sigma) \inf_{\substack{C \subseteq [p] \\ |C|=k}} \inf_{\mathbf{u} \in V_C} \frac{\|\mathbf{Z}\mathbf{u}\|}{\sqrt{n}\|\mathbf{u}\|},$$

where $\mathbf{Z} = \mathbf{X}\Sigma^{-1/2} \in \mathbb{R}^{n \times p}$ is a matrix with iid $\mathcal{N}(0, 1)$ entries, and $V_C = \{\Sigma^{1/2}\mathbf{b} : \mathbf{b} \in \mathbb{R}^p, \text{supp}(\mathbf{b}) \subseteq C \subseteq [p], |C| = k\}$. The linear subspace V_C has dimension k and we consider a matrix $\mathbf{Q} \in \mathbb{R}^{p \times k}$ such

that its columns form an orthonormal basis of V_C . By construction, QQ^\top is the orthogonal projection onto V_C and $Q^\top Q = I_k$. For any $\mathbf{u} \in V_C$, let $\mathbf{v} = Q^\top \mathbf{u}$, so that

$$\frac{\|\mathbf{Z}\mathbf{u}\|}{\sqrt{n}\|\mathbf{u}\|} = \frac{\|\mathbf{Z}Q\mathbf{v}\|}{\sqrt{n}\|\mathbf{v}\|} \geq \frac{\phi_{\min}(\mathbf{Z}Q)}{\sqrt{n}},$$

where the matrix $\mathbf{Z}Q \in \mathbb{R}^{n \times k}$ has iid $\mathcal{N}(0, 1)$ entries since $Q^\top Q = I_k$. Using (Davidson and Szarek, 2001, Theorem II.13) on the smallest singular value of a Gaussian matrix, we have for every $t \geq 0$

$$\mathbb{P}(\phi_{\min}(\mathbf{Z}Q) \geq \sqrt{n} - \sqrt{k} - t) \geq 1 - e^{-t^2/2}. \quad (27)$$

We take $t = \frac{\sqrt{n}}{4}$, and choose k such that $\sqrt{k} \leq \frac{\sqrt{n}}{4}$, so that $\sqrt{n} - \sqrt{k} - t \geq \frac{\sqrt{n}}{2}$. Taking all the $\binom{p}{k}$ combinations of support set C such that $C \subseteq [p]$ and $|C| = k$, we obtain

$$\mathbb{P}\left(\inf_{\mathbf{b} \in \mathbb{R}^p: \|\mathbf{b}\|_0 \leq k} \frac{\|\mathbf{X}\mathbf{b}\|}{\sqrt{n}\|\mathbf{b}\|} \geq \phi_{\min}^{1/2}(\Sigma)/2\right) \geq 1 - \binom{p}{k} e^{-n/32} \geq 1 - e^{k \log(ep/k) - n/32}.$$

We take $k = c^*(\gamma)n$ with $c^*(\gamma) = \sup_{c \in [0, \frac{1}{16} \wedge \gamma]} \{c \log(ep/c) \leq 1/64\}$, then

$$\mathbb{P}\left(\inf_{\mathbf{b} \in \mathbb{R}^p: \|\mathbf{b}\|_0 \leq k} \frac{\|\mathbf{X}\mathbf{b}\|}{\sqrt{n}\|\mathbf{b}\|} \geq \phi_{\min}^{1/2}(\Sigma)/2\right) \geq 1 - e^{-n/64}.$$

With $\delta = \phi_{\min}^{1/2}(\Sigma)/2$, and $k = c^*(\gamma)n$, then condition (26) holds with probability at least $1 - e^{-n/64}$. In other words, letting

$$A_0 = \left\{ \mathbf{b} \in \mathbb{R}^p: \inf_{\|\mathbf{b}\|_0 \leq c^*(\gamma)n} \frac{\|\mathbf{X}\mathbf{b}\|}{\sqrt{n}\|\mathbf{b}\|} \geq \phi_{\min}^{1/2}(\Sigma)/2 \right\},$$

then $\mathbb{P}(A_0) \geq 1 - e^{-n/64}$. Now we are ready to bound κ^2 in Lemma A.1 using Lemma A.2. To this end, we first derive a cone condition for $\mathbf{B} - \mathbf{B}^*$ using $\|\mathbf{B}\|_{2,1} \leq \|\mathbf{B}^*\|_{2,1}$. By rearranging the terms and Cauchy-Schwarz inequality, we have

$$\sum_{j \in \mathcal{S}^c} \|\mathbf{B}^\top \mathbf{e}_j\| \leq \sum_{j \in \mathcal{S}} (\|\mathbf{B}^{*\top} \mathbf{e}_j\| - \|\mathbf{B}^\top \mathbf{e}_j\|) \leq \sum_{j \in \mathcal{S}} \|(\mathbf{B} - \mathbf{B}^*)^\top \mathbf{e}_j\| \leq \sqrt{|\mathcal{S}|} \|\mathbf{B} - \mathbf{B}^*\|_F.$$

Since $\sum_{j \in \mathcal{S}^c} \|(\mathbf{B} - \mathbf{B}^*)^\top \mathbf{e}_j\| = \sum_{j \in \mathcal{S}^c} \|\mathbf{B}^\top \mathbf{e}_j\| \leq \sqrt{|\mathcal{S}|} \|\mathbf{B} - \mathbf{B}^*\|_F$, we obtain the cone condition

$$\sum_{j=1}^p \|(\mathbf{B} - \mathbf{B}^*)^\top \mathbf{e}_j\| \leq 2\sqrt{|\mathcal{S}|} \|\mathbf{B} - \mathbf{B}^*\|_F. \quad (28)$$

Let $A_1 = \left\{ \max_{j \in [p]} \frac{\|\mathbf{X}\mathbf{e}_j\|^2}{n} \leq 1.1 \right\}$, then under the assumption that $\Sigma_{jj} = 1$ for all $j \in [p]$, we have $\mathbb{P}(A_1) \geq 1 - e^{-c(\gamma)n}$ by (Laurent and Massart, 2000, Lemma 1). Now we apply Lemma A.2 to $\mathbf{A} = \mathbf{B} - \mathbf{B}^*$, using the cone (28) and $\sum_j \mu_j = 1$, we have on the event A_1 ,

$$\frac{\|\mathbf{X}(\mathbf{B} - \mathbf{B}^*)\|_F^2}{n} \geq \delta^2 \left(1 - \frac{1}{k}\right) \|\mathbf{B} - \mathbf{B}^*\|_F^2 - \frac{4.4|\mathcal{S}|}{k} \|\mathbf{B} - \mathbf{B}^*\|_F^2.$$

Since $1 - \frac{1}{k} = 1 - \frac{1}{c^*(\gamma)n} \geq 1/2$ for n large enough, with $|\mathcal{S}| \leq \frac{k\delta^2}{17.6} \leq \frac{k\delta^2}{16}$, we have

$$\frac{\|\mathbf{X}(\mathbf{B} - \mathbf{B}^*)\|_{\mathbb{F}}^2}{n} \geq \frac{\delta^2}{4} \|\mathbf{B} - \mathbf{B}^*\|_{\mathbb{F}}^2, \quad (29)$$

and thus $\kappa \geq \frac{\delta}{2}$ with $\delta = \phi_{\min}^{1/2}(\Sigma)/2$.

Now we are ready to bound $\|\widehat{\mathbf{B}}\|_0$ using Lemma A.1. Let $A_2 = \{\|\mathbf{X}^\top \mathbf{X}\|_{\text{op}} \leq \|\Sigma\|_{\text{op}}(2 + \sqrt{\gamma})^2 n\}$, $A_3 = \{\|\mathbf{E}\|_{\mathbb{F}}^2 \leq 5 \text{Tr}(\mathbf{S})n\}$, then $\mathbb{P}(A_2) \geq 1 - \exp(-n/2)$ by (Davidson and Szarek, 2001, Theorem II.13), and $\mathbb{P}(A_3) \geq 1 - \exp(-n)$ by (Laurent and Massart, 2000, Lemma 1). On the event $A_0 \cap A_1 \cap A_2 \cap A_3$, we have

$$\begin{aligned} \|\widehat{\mathbf{B}}\|_0 &\leq \|\mathbf{X}\|_{\text{op}}^2 \max \left\{ \frac{2\|\mathbf{E}\|_{\mathbb{F}}^2}{n^2 \lambda^2}, \frac{4\|\mathbf{B}^*\|_0}{n \kappa^2} \right\} \\ &\leq \|\Sigma\|_{\text{op}}(2 + \sqrt{\gamma})^2 \max \left\{ \frac{10 \text{Tr}(\mathbf{S})}{\lambda^2}, \frac{4\|\mathbf{B}^*\|_0}{\kappa^2} \right\} \\ &\leq \|\Sigma\|_{\text{op}}(2 + \sqrt{\gamma})^2 \max \left\{ \frac{10n}{\mu^2}, \frac{64\|\mathbf{B}^*\|_0}{\phi_{\min}(\Sigma)} \right\}, \end{aligned}$$

where the last inequality uses $\lambda = \mu\sqrt{\text{Tr}(\mathbf{S})/n}$, and $\kappa \geq \delta/2 = \phi_{\min}^{1/2}(\Sigma)/4$.

To guarantee $\|\widehat{\mathbf{B}}\|_0 \leq n/3$, we need (1) $\|\Sigma\|_{\text{op}}(2 + \sqrt{\gamma})^2 \frac{10}{\mu^2} \leq 1/3$, that is, $\mu \geq \mu^*(\gamma, \Sigma) := (30\|\Sigma\|_{\text{op}})^{1/2}(2 + \sqrt{\gamma})$, and (2) $\|\Sigma\|_{\text{op}}(2 + \sqrt{\gamma})^2 \frac{64\|\mathbf{B}^*\|_0}{\phi_{\min}(\Sigma)} \leq n/3$, that is, $\|\mathbf{B}^*\|_0 \leq \frac{\phi_{\min}(\Sigma)}{192\|\Sigma\|_{\text{op}}(2 + \sqrt{\gamma})^2} n$. Recall that we also need $\|\mathbf{B}^*\|_0 \leq \frac{k\delta^2}{16} \leq \frac{c^*(\gamma)\phi_{\min}(\Sigma)}{64} n := c^*(\gamma, \Sigma)n$ to ensure (29). So under the Assumption 4(iii), on the event $A_0 \cap A_1 \cap A_2 \cap A_3$, we have $\|\widehat{\mathbf{B}}\|_0 \leq n/3$. Using the union bound for the events A_0, A_1, A_2, A_3 , we have

$$\mathbb{P}(\|\widehat{\mathbf{B}}\|_0 \leq n/3) \geq 1 - \exp(-c(\gamma)n) > 1 - \frac{1}{T}$$

if $T \leq e^{\sqrt{n}}$. ■

To complete the proof of Proposition 3.3, we now provide proofs for Lemmas A.1 and A.2.

Proof of Lemma A.1. The KKT conditions for $\widehat{\mathbf{B}}$ read $\mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\widehat{\mathbf{B}}) = n\lambda\partial\|\widehat{\mathbf{B}}\|_{2,1}$. Let $\hat{\mathcal{S}}$ be the support of $\widehat{\mathbf{B}}$, then for each $j \in \hat{\mathcal{S}}$, we have $\mathbf{e}_j^\top \mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\widehat{\mathbf{B}}) = n\lambda \frac{\mathbf{e}_j^\top \widehat{\mathbf{B}}}{\|\widehat{\mathbf{B}}\|_{2,1}}$. Taking squared norm on both sides and summing over $j \in \hat{\mathcal{S}}$ gives

$$(n\lambda)^2 |\hat{\mathcal{S}}| = \sum_{j \in \hat{\mathcal{S}}} \|\mathbf{e}_j^\top \mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\widehat{\mathbf{B}})\|^2 \leq \|\mathbf{Y} - \mathbf{X}\widehat{\mathbf{B}}\|_{\mathbb{F}}^2 \|\mathbf{X}^\top \mathbf{X}\|_{\text{op}}. \quad (30)$$

Now we bound $\|\mathbf{Y} - \mathbf{X}\widehat{\mathbf{B}}\|_{\mathbb{F}}^2$. Multiplying the KKT conditions by $\widehat{\mathbf{B}} - \mathbf{B}^*$ gives

$$\|\mathbf{X}(\widehat{\mathbf{B}} - \mathbf{B}^*)\|_{\mathbb{F}}^2 + \|\mathbf{Y} - \mathbf{X}\widehat{\mathbf{B}}\|_{\mathbb{F}}^2 \leq \|\mathbf{E}\|_{\mathbb{F}}^2 + 2n\lambda(\|\mathbf{B}^*\|_{2,1} - \|\widehat{\mathbf{B}}\|_{2,1}).$$

We bound the RHS by comparing the two terms. If $\|\mathbf{E}\|_{\mathbb{F}}^2 \geq 2n\lambda(\|\mathbf{B}^*\|_{2,1} - \|\widehat{\mathbf{B}}\|_{2,1})$, then

$$\|\mathbf{Y} - \mathbf{X}\widehat{\mathbf{B}}\|_{\mathbb{F}}^2 \leq (RHS) \leq 2\|\mathbf{E}\|_{\mathbb{F}}^2.$$

Otherwise $\|E\|_F^2 \leq 2n\lambda(\|\mathbf{B}^*\|_{2,1} - \|\widehat{\mathbf{B}}\|_{2,1})$, then we have

$$\begin{aligned} \|\mathbf{X}(\widehat{\mathbf{B}} - \mathbf{B}^*)\|_F^2 + \|\mathbf{Y} - \mathbf{X}\widehat{\mathbf{B}}\|_F^2 &\leq 4n\lambda(\|\mathbf{B}^*\|_{2,1} - \|\widehat{\mathbf{B}}\|_{2,1}) \\ &\leq 4n\lambda\sqrt{|\mathcal{S}|}\|\mathbf{B}^* - \widehat{\mathbf{B}}\|_F \\ &\leq 4\sqrt{n}\lambda\sqrt{|\mathcal{S}|}\kappa^{-1}\|\mathbf{X}(\widehat{\mathbf{B}} - \mathbf{B}^*)\|_F. \end{aligned}$$

where \mathcal{S} is the support of \mathbf{B}^* . Using $4ab \leq 4a^2 + b^2$ with $a = \sqrt{n}\lambda\sqrt{|\mathcal{S}|}\kappa^{-1}$ and $b = \|\mathbf{X}(\widehat{\mathbf{B}} - \mathbf{B}^*)\|_F$, then the term b^2 cancels out in both sides. We have

$$\|\mathbf{Y} - \mathbf{X}\widehat{\mathbf{B}}\|_F^2 \leq 4n\lambda^2|\mathcal{S}|\kappa^{-2}.$$

Plugging the maximum of the previous two bounds for $\|\mathbf{Y} - \mathbf{X}\widehat{\mathbf{B}}\|_F^2$ into (30) gives the desired inequality. \blacksquare

Proof of Lemma A.2. Replacing \mathbf{A} by $\mathbf{A}/\sum_{j=1}^P \|\mathbf{A}^\top \mathbf{e}_j\|$ if necessary, we assume without loss of generality that $\sum_{j=1}^P \|\mathbf{A}^\top \mathbf{e}_j\| = 1$ and we also assume for simplicity that $\forall j \in [P], \mathbf{A}^\top \mathbf{e}_j \neq 0$. Define a discrete random matrix $\widetilde{\mathbf{A}}$ by $\mathbb{P}(\widetilde{\mathbf{A}} = \mathbf{e}_j \mathbf{e}_j^\top \mathbf{A} / \mu_j) = \mu_j$, where $\mu_j = \|\mathbf{A}^\top \mathbf{e}_j\|$. Then we have $\mathbb{E}\widetilde{\mathbf{A}} = \mathbf{A}$. Let $\widetilde{\mathbf{A}}_1, \dots, \widetilde{\mathbf{A}}_k$ be iid copies of $\widetilde{\mathbf{A}}$ and set $\bar{\mathbf{A}} = \frac{1}{k} \sum_{m=1}^k \widetilde{\mathbf{A}}_m$, then $\bar{\mathbf{A}}$ is k -row sparse by construction and condition (26) gives $\|\mathbf{X}\bar{\mathbf{A}}\|_F^2/n \geq \delta^2\|\bar{\mathbf{A}}\|_F^2$. Hence

$$\mathbb{E}[\|\mathbf{X}\bar{\mathbf{A}}\|_F^2/n] \geq \delta^2\mathbb{E}[\|\bar{\mathbf{A}}\|_F^2]. \quad (31)$$

Note that all expectations are conditionally on \mathbf{X} and the randomness stems from $\widetilde{\mathbf{A}}_1, \dots, \widetilde{\mathbf{A}}_k$. Now we calculate the expectations (conditionally on \mathbf{X}) by expanding the squares. For $u \neq v$, we have $\mathbb{E}\langle \mathbf{X}\widetilde{\mathbf{A}}_u, \mathbf{X}\widetilde{\mathbf{A}}_v \rangle = \|\mathbf{X}\mathbf{A}\|_F^2$ by independence. For $u = v$, using $\mu_j = \|\mathbf{A}^\top \mathbf{e}_j\|$, we have

$$\mathbb{E}[\|\mathbf{X}\widetilde{\mathbf{A}}_u\|_F^2/n] = \sum_j \mu_j \|\mathbf{X}\mathbf{e}_j \mathbf{e}_j^\top \mathbf{A}\|_F^2 / (\mu_j^2 n) = \sum_j \mu_j \|\mathbf{X}\mathbf{e}_j\|^2 / n.$$

Therefore,

$$\mathbb{E}[\|\mathbf{X}\bar{\mathbf{A}}\|_F^2/n] = \left(1 - \frac{1}{k}\right) \|\mathbf{X}\mathbf{A}\|_F^2/n + \frac{1}{k} \sum_j \mu_j \|\mathbf{X}\mathbf{e}_j\|^2 / n.$$

The same argument gives

$$\mathbb{E}[\|\bar{\mathbf{A}}\|_F^2/n] = \left(1 - \frac{1}{k}\right) \|\mathbf{A}\|_F^2 + \frac{1}{k} \sum_j \mu_j = \left(1 - \frac{1}{k}\right) \|\mathbf{A}\|_F^2 + \frac{1}{k}.$$

Plugging the above two equalities into (31) gives the desired bound. \blacksquare

Appendix B: Proof of main results

We first present two lemmas and three propositions in Appendix B.1. Based on these technical results, we prove the main results of the paper in Appendix B.2. The proofs of technical results in Appendix B.1 are deferred to the supplementary material Tan, Romon and Bellec (2023).

B.1. Preliminary results

The following Lemmas B.1 and B.2 provide operator norm bounds for $\mathbf{I}_T - \widehat{\mathbf{A}}/n$ and $(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1}$ under suitable conditions.

Lemma B.1. *Suppose that Assumption 2 holds. If $\tau > 0$ in (6) with $\tau' = \tau / \|\Sigma\|_{\text{op}}$, then*

- (i) $\|\mathbf{I}_T - \widehat{\mathbf{A}}/n\|_{\text{op}} \leq 1$.
- (ii) *In the event $\{\|\mathbf{X}\Sigma^{-\frac{1}{2}}\|_{\text{op}} < 2\sqrt{n} + \sqrt{p}\}$, we have $\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1}\|_{\text{op}} \leq 1 + (\tau')^{-1}(2 + \sqrt{p/n})^2$. Furthermore, $\mathbb{E}[\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1}\|_{\text{op}}] \leq 1 + (\tau')^{-1}[(1 + \sqrt{p/n})^2 + n^{-1}]$.*

Lemma B.2. *Suppose that Assumption 2 holds. If $\tau = 0$ in (6), then*

- (i) *In the event U_1 , we have $\|\mathbf{I}_T - \widehat{\mathbf{A}}/n\|_{\text{op}} \leq 1$.*
- (ii) *In the event U_1 , $\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1}\|_{\text{op}} \leq C(c)$. Hence, $\mathbb{E}[I(U_1)\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1}\|_{\text{op}}] \leq C(c)$.*

The proof of main results rely on the following three propositions, which are derived from non-trivial applications of second-order Stein's identity (Bellec and Zhang, 2021).

Proposition B.3. *Suppose that Assumption 1 holds.*

Let $\mathbf{Q}_1 = \frac{\frac{1}{n}(\mathbf{F}^\top \mathbf{F} + \mathbf{H}^\top \mathbf{Z}^\top \mathbf{F} - \mathbf{S}(n\mathbf{I}_T - \widehat{\mathbf{A}}))}{\|S^{\frac{1}{2}}\|_{\text{F}}(\|\mathbf{F}\|_{\text{F}}^2/n + \text{Tr}(\mathbf{S}))^{\frac{1}{2}}n^{-\frac{1}{2}}}$, then $\mathbb{E}[\|\mathbf{Q}_1\|_{\text{F}}^2] \leq 4$.

Proposition B.4. *Suppose that Assumptions 2, 3 and 4 hold. Let*

$$\mathbf{Q}_2 = \frac{\frac{1}{n^2}(\mathbf{F}^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{F} - \mathbf{F}^\top \mathbf{F}(p\mathbf{I}_T - \widehat{\mathbf{A}}) + (n\mathbf{I}_T - \widehat{\mathbf{A}})\mathbf{H}^\top \mathbf{Z}^\top \mathbf{F})}{(\|\mathbf{H}\|_{\text{F}}^2 + \|\mathbf{F}\|_{\text{F}}^2/n)n^{-\frac{1}{2}}},$$

then $\mathbb{E}[\|\mathbf{Q}_2\|_{\text{F}}^2] \leq C(\tau')(T \wedge (1 + \frac{p}{n}))(1 + \frac{p}{n})$ under Assumption 4(i), and $\mathbb{E}[I(\Omega)\|\mathbf{Q}_2\|_{\text{F}}^2] \leq C(\gamma, c)$ under Assumption 4(ii) for some set Ω with $\mathbb{P}(\Omega) \rightarrow 1$.

Proposition B.5. *Suppose that Assumptions 2, 3 and 4 hold. Let $\Xi = (n\mathbf{I}_T - \widehat{\mathbf{A}})\mathbf{H}^\top \mathbf{Z}^\top \mathbf{F}$, and*

$$\mathbf{Q}_3 = \frac{\frac{1}{n^2}(p\mathbf{F}^\top \mathbf{F} - \mathbf{F}^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{F} - (n\mathbf{I}_T - \widehat{\mathbf{A}})\mathbf{H}^\top \mathbf{H}(n\mathbf{I}_T - \widehat{\mathbf{A}}) - \Xi - \Xi^\top)}{(\|\mathbf{H}\|_{\text{F}}^2 + \|\mathbf{F}\|_{\text{F}}^2/n)n^{-\frac{1}{2}}},$$

then $\mathbb{E}[\|\mathbf{Q}_3\|_{\text{F}}] \leq C(\gamma, \tau')$ under Assumption 4(i), and $\mathbb{E}[I(\Omega)\|\mathbf{Q}_3\|_{\text{F}}] \leq C(\gamma, c)$ under Assumption 4(ii) for some set Ω with $\mathbb{P}(\Omega) \rightarrow 1$.

B.2. Proofs of main results

With the preliminary results in previous section, we are ready to present the proofs of main results.

Proof of Theorem 3.4. Recall definition of $\widehat{\mathbf{S}}$ in Definition 2.2, and let $\mathbf{Q}_1, \mathbf{Q}_2$ be defined as in Propositions B.3 and B.4. With $\mathbf{Z} = \mathbf{X}\Sigma^{-1/2}$, we obtain

$$n^2 \left[\mathbf{Q}_2(\|\mathbf{H}\|_{\text{F}}^2 + \|\mathbf{F}\|_{\text{F}}^2/n)n^{-\frac{1}{2}} - n^{-1}(n\mathbf{I}_T - \widehat{\mathbf{A}})\mathbf{Q}_1(\|S^{\frac{1}{2}}\|_{\text{F}}(\|\mathbf{F}\|_{\text{F}}^2/n + \text{Tr}(\mathbf{S}))^{\frac{1}{2}}n^{-\frac{1}{2}}) \right]$$

$$\begin{aligned}
&= (\mathbf{F}^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{F} - \mathbf{F}^\top \mathbf{F} (p \mathbf{I}_T - \widehat{\mathbf{A}}) + (n \mathbf{I}_T - \widehat{\mathbf{A}}) \mathbf{H}^\top \mathbf{Z}^\top \mathbf{F}) \\
&\quad - [(n \mathbf{I}_T - \widehat{\mathbf{A}}) (\mathbf{F}^\top \mathbf{F} + \mathbf{H}^\top \mathbf{Z}^\top \mathbf{F} - \mathbf{S} (n \mathbf{I}_T - \widehat{\mathbf{A}}))] \\
&= (\mathbf{F}^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{F} - \mathbf{F}^\top \mathbf{F} (p \mathbf{I}_T - \widehat{\mathbf{A}})) - (n \mathbf{I}_T - \widehat{\mathbf{A}}) (\mathbf{F}^\top \mathbf{F} - \mathbf{S} (n \mathbf{I}_T - \widehat{\mathbf{A}})) \\
&= \mathbf{F}^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{F} + \mathbf{F}^\top \mathbf{F} \widehat{\mathbf{A}} + \widehat{\mathbf{A}} \mathbf{F}^\top \mathbf{F} - (n + p) \mathbf{F}^\top \mathbf{F} + (n \mathbf{I}_T - \widehat{\mathbf{A}}) \mathbf{S} (n \mathbf{I}_T - \widehat{\mathbf{A}}) \\
&= (n \mathbf{I}_T - \widehat{\mathbf{A}}) \mathbf{S} (n \mathbf{I}_T - \widehat{\mathbf{A}}) + \mathbf{F}^\top \mathbf{F} \widehat{\mathbf{A}} + \widehat{\mathbf{A}} \mathbf{F}^\top \mathbf{F} - \mathbf{F}^\top ((n + p) \mathbf{I}_T - \mathbf{Z} \mathbf{Z}^\top) \mathbf{F} \\
&= (n \mathbf{I}_T - \widehat{\mathbf{A}}) \mathbf{S} (n \mathbf{I}_T - \widehat{\mathbf{A}}) - (n \mathbf{I}_T - \widehat{\mathbf{A}}) \widehat{\mathbf{S}} (n \mathbf{I}_T - \widehat{\mathbf{A}}) \\
&= (n \mathbf{I}_T - \widehat{\mathbf{A}}) (\mathbf{S} - \widehat{\mathbf{S}}) (n \mathbf{I}_T - \widehat{\mathbf{A}}).
\end{aligned}$$

Therefore, by the triangle inequality and $\|\mathbf{I}_T - \widehat{\mathbf{A}}/n\|_{\text{op}} \leq 1$ in Lemmas B.1 and B.2,

$$\begin{aligned}
&\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)(\mathbf{S} - \widehat{\mathbf{S}})(\mathbf{I}_T - \widehat{\mathbf{A}}/n)\|_{\text{F}} \\
&\leq \|\mathbf{Q}_2\|_{\text{F}} n^{-\frac{1}{2}} (\|\mathbf{H}\|_{\text{F}}^2 + \|\mathbf{F}\|_{\text{F}}^2/n) + \|\mathbf{I}_T - \widehat{\mathbf{A}}/n\|_{\text{op}} \|\mathbf{Q}_1\|_{\text{F}} n^{-\frac{1}{2}} \|\mathbf{S}^{\frac{1}{2}}\|_{\text{F}} (\|\mathbf{F}\|_{\text{F}}^2/n + \text{Tr}(\mathbf{S}))^{\frac{1}{2}} \\
&\leq \|\mathbf{Q}_2\|_{\text{F}} n^{-\frac{1}{2}} (\|\mathbf{H}\|_{\text{F}}^2 + \|\mathbf{F}\|_{\text{F}}^2/n) + \|\mathbf{Q}_1\|_{\text{F}} n^{-\frac{1}{2}} \|\mathbf{S}^{\frac{1}{2}}\|_{\text{F}} (\|\mathbf{F}\|_{\text{F}}^2/n + \text{Tr}(\mathbf{S}))^{\frac{1}{2}} \\
&\leq \|\mathbf{Q}_2\|_{\text{F}} n^{-\frac{1}{2}} (\|\mathbf{H}\|_{\text{F}}^2 + \|\mathbf{F}\|_{\text{F}}^2/n) + \|\mathbf{Q}_1\|_{\text{F}} n^{-\frac{1}{2}} \frac{1}{2} [\text{Tr}(\mathbf{S}) + (\|\mathbf{F}\|_{\text{F}}^2/n + \text{Tr}(\mathbf{S}))] \\
&\leq (\|\mathbf{Q}_2\|_{\text{F}} + \|\mathbf{Q}_1\|_{\text{F}}) n^{-\frac{1}{2}} (\|\mathbf{H}\|_{\text{F}}^2 + \|\mathbf{F}\|_{\text{F}}^2/n + \text{Tr}(\mathbf{S})).
\end{aligned}$$

Therefore,

$$\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)(\mathbf{S} - \widehat{\mathbf{S}})(\mathbf{I}_T - \widehat{\mathbf{A}}/n)\|_{\text{F}} \leq \Theta_1 n^{-\frac{1}{2}} (\|\mathbf{H}\|_{\text{F}}^2 + \|\mathbf{F}\|_{\text{F}}^2/n + \text{Tr}(\mathbf{S})),$$

where $\Theta_1 = \|\mathbf{Q}_1\|_{\text{F}} + \|\mathbf{Q}_2\|_{\text{F}}$. Note that we have $\mathbb{E}[\|\mathbf{Q}_1\|_{\text{F}}^2] \leq 4$ from Proposition B.3. By Proposition B.4, we have

(1) under Assumption 4(i), $\mathbb{E}[\|\mathbf{Q}_2\|_{\text{F}}^2] \leq C(\tau')(T \wedge (1 + \frac{p}{n}))(1 + \frac{p}{n})$. Hence

$$\begin{aligned}
\mathbb{E}[\Theta_1^2] &\leq 2\mathbb{E}[\|\mathbf{Q}_1\|_{\text{F}}^2 + \|\mathbf{Q}_2\|_{\text{F}}^2] \leq 2[4 + C(\tau')(T \wedge (1 + \frac{p}{n}))(1 + \frac{p}{n})] \\
&\leq C(\tau')(T \wedge (1 + \frac{p}{n}))(1 + \frac{p}{n}).
\end{aligned}$$

(2) under Assumption 4(ii), $\mathbb{E}[I(\Omega)\|\mathbf{Q}_2\|_{\text{F}}^2] \leq C(\gamma, c)$ with $\mathbb{P}(\Omega) \rightarrow 1$. Thus,

$$\mathbb{E}[I(\Omega)\Theta_1^2] \leq 2\mathbb{E}[\|\mathbf{Q}_1\|_{\text{F}}^2 + I(\Omega)\|\mathbf{Q}_2\|_{\text{F}}^2] \leq C(\gamma, c).$$

■

Proof of Theorem 3.9. From the definitions of $\mathbf{Q}_1, \mathbf{Q}_2, \mathbf{Q}_3$ in Propositions B.3, B.4 and B.5, we have

$$\begin{aligned}
&(\mathbf{Q}_2^\top + \mathbf{Q}_3) n^{-\frac{1}{2}} (\|\mathbf{H}\|_{\text{F}}^2 + \|\mathbf{F}\|_{\text{F}}^2/n) + (\mathbf{I}_T - \widehat{\mathbf{A}}/n) \mathbf{Q}_1 n^{-\frac{1}{2}} \|\mathbf{S}^{\frac{1}{2}}\|_{\text{F}} (\|\mathbf{F}\|_{\text{F}}^2/n + \text{Tr}(\mathbf{S}))^{\frac{1}{2}} \\
&= \frac{1}{n} \mathbf{F}^\top \mathbf{F} - (\mathbf{I}_T - \widehat{\mathbf{A}}/n) (\mathbf{H}^\top \mathbf{H} + \mathbf{S}) (\mathbf{I}_T - \widehat{\mathbf{A}}/n) \\
&= (\mathbf{I}_T - \widehat{\mathbf{A}}/n) [n^{-1} (\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1} \mathbf{F}^\top \mathbf{F} (\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1} - (\mathbf{H}^\top \mathbf{H} + \mathbf{S})] (\mathbf{I}_T - \widehat{\mathbf{A}}/n)
\end{aligned}$$

$$= (\mathbf{I}_T - \widehat{\mathbf{A}}/n)(\widehat{\mathbf{R}} - \mathbf{R})(\mathbf{I}_T - \widehat{\mathbf{A}}/n),$$

where $\widehat{\mathbf{R}} \stackrel{\text{def}}{=} n^{-1}(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1} \mathbf{F}^\top \mathbf{F}(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1}$, and $\mathbf{R} \stackrel{\text{def}}{=} \mathbf{H}^\top \mathbf{H} + \mathbf{S}$.

Therefore, by the triangle inequality and $\|\mathbf{I}_T - \widehat{\mathbf{A}}/n\|_{\text{op}} \leq 1$ from Lemmas B.1 and B.2,

$$\begin{aligned} & \|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)(\widehat{\mathbf{R}} - \mathbf{R})(\mathbf{I}_T - \widehat{\mathbf{A}}/n)\|_{\text{F}} \\ & \leq (\|\mathbf{Q}_2\|_{\text{F}} + \|\mathbf{Q}_3\|_{\text{F}})n^{-\frac{1}{2}}(\|\mathbf{H}\|_{\text{F}}^2 + \|\mathbf{F}\|_{\text{F}}^2/n) + \|\mathbf{Q}_1\|_{\text{F}}n^{-\frac{1}{2}}\|\mathbf{S}\|_{\text{F}}^{\frac{1}{2}}(\|\mathbf{F}\|_{\text{F}}^2/n + \text{Tr}(\mathbf{S}))^{\frac{1}{2}} \\ & \leq (\|\mathbf{Q}_2\|_{\text{F}} + \|\mathbf{Q}_3\|_{\text{F}} + \|\mathbf{Q}_1\|_{\text{F}})n^{-\frac{1}{2}}(\|\mathbf{F}\|_{\text{F}}^2/n + \|\mathbf{H}\|_{\text{F}}^2 + \text{Tr}(\mathbf{S})) \\ & = \Theta_2 n^{-\frac{1}{2}}(\|\mathbf{F}\|_{\text{F}}^2/n + \|\mathbf{H}\|_{\text{F}}^2 + \text{Tr}(\mathbf{S})), \end{aligned}$$

where $\Theta_2 = \|\mathbf{Q}_1\|_{\text{F}} + \|\mathbf{Q}_2\|_{\text{F}} + \|\mathbf{Q}_3\|_{\text{F}}$. By Propositions B.3, B.4 and B.5, we obtain $\mathbb{E}[\Theta_2] \leq C(\gamma, \tau')$ under Assumption 4(i) and $\mathbb{E}[I(\Omega)\Theta_2] \leq C(\gamma, c)$ with $P(\Omega) \rightarrow 1$ under Assumption 4(ii).

Furthermore, since $\Theta_2 = O_P(1)$, and $\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1}\|_{\text{op}} = O_P(1)$ from Lemmas B.1 and B.2,

$$\begin{aligned} \|\widehat{\mathbf{R}} - \mathbf{R}\|_{\text{F}} & \leq \|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1}\|_{\text{op}}^2 \Theta_2 n^{-\frac{1}{2}}(\|\mathbf{F}\|_{\text{F}}^2/n + \|\mathbf{H}\|_{\text{F}}^2 + \text{Tr}(\mathbf{S})) \\ & = O_P(n^{-\frac{1}{2}})(\|\mathbf{F}\|_{\text{F}}^2/n + \|\mathbf{H}\|_{\text{F}}^2 + \text{Tr}(\mathbf{S})). \end{aligned}$$

Since $\frac{1}{n}\mathbf{F}^\top \mathbf{F} = (\mathbf{I}_T - \widehat{\mathbf{A}}/n)\widehat{\mathbf{R}}(\mathbf{I}_T - \widehat{\mathbf{A}}/n)$, taking trace of both sides gives $\frac{1}{n}\|\mathbf{F}\|_{\text{F}}^2 \leq \|\widehat{\mathbf{R}}\|_*$ thanks to $\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)\|_{\text{op}} \leq 1$. Note that $\|\mathbf{R}\|_* = \|\mathbf{H}\|_{\text{F}}^2 + \text{Tr}(\mathbf{S})$ by definition of \mathbf{R} , we obtain

$$\|\widehat{\mathbf{R}} - \mathbf{R}\|_{\text{F}} \leq O_P(n^{-\frac{1}{2}})(\|\widehat{\mathbf{R}}\|_* + \|\mathbf{R}\|_*). \quad (32)$$

Since $\widehat{\mathbf{R}}$ and \mathbf{R} are both $T \times T$ positive semi-definite matrices, whose ranks are at most T ,

$$\begin{aligned} |\|\widehat{\mathbf{R}}\|_* - \|\mathbf{R}\|_*| & \leq \|\widehat{\mathbf{R}} - \mathbf{R}\|_* \leq \sqrt{T}\|\widehat{\mathbf{R}} - \mathbf{R}\|_{\text{F}} \\ & \leq O_P((T/n)^{\frac{1}{2}})(\|\widehat{\mathbf{R}}\|_* + \|\mathbf{R}\|_*) = o_P(1)(\|\widehat{\mathbf{R}}\|_* + \|\mathbf{R}\|_*), \end{aligned}$$

thanks to $T = o(n)$. That is,

$$\frac{|\|\widehat{\mathbf{R}}\|_* - \|\mathbf{R}\|_*|}{\|\widehat{\mathbf{R}}\|_* + \|\mathbf{R}\|_*} \leq O_P((T/n)^{\frac{1}{2}}),$$

which implies $\frac{\|\mathbf{R}\|_*}{\|\widehat{\mathbf{R}}\|_*} - 1 = O_P((T/n)^{\frac{1}{2}})$, i.e.,

$$\frac{\text{Tr}(\mathbf{S}) + \|\mathbf{H}\|_{\text{F}}^2}{\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1} \mathbf{F}^\top\|_{\text{F}}^2/n} - 1 = O_P((T/n)^{\frac{1}{2}}) = o_P(1).$$

■

Proof of Theorem 3.5. This proof is based on results of Theorems 3.4 and 3.9. We begin with the result of Theorem 3.9,

$$\frac{\text{Tr}(\mathbf{S}) + \|\mathbf{H}\|_{\text{F}}^2}{\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1} \mathbf{F}^\top\|_{\text{F}}^2/n} \xrightarrow{P} 1.$$

In other words,

$$\text{Tr}(\mathbf{S}) + \|\mathbf{H}\|_{\mathbb{F}}^2 = (1 + o_P(1)) \|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1} \mathbf{F}^\top\|_{\mathbb{F}}^2/n.$$

Thus, the upper bound in Theorem 3.4 can be bounded from above as follows

$$\begin{aligned} & \|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)(\widehat{\mathbf{S}} - \mathbf{S})(\mathbf{I}_T - \widehat{\mathbf{A}}/n)\|_{\mathbb{F}} \\ & \leq \Theta_1 n^{-\frac{1}{2}} (\|\mathbf{F}\|_{\mathbb{F}}^2/n + \|\mathbf{H}\|_{\mathbb{F}}^2 + \|\mathbf{S}^{\frac{1}{2}}\|_{\mathbb{F}}^2) \\ & \leq \Theta_1 n^{-\frac{1}{2}} (\|\mathbf{F}\|_{\mathbb{F}}^2/n + (1 + o_P(1)) \|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1} \mathbf{F}^\top\|_{\mathbb{F}}^2/n) \\ & \leq \Theta_1 n^{-\frac{1}{2}} (1 + (1 + o_P(1)) \|(\mathbf{I}_T - \widehat{\mathbf{A}})^{-1}\|_{\text{op}}^2) \|\mathbf{F}\|_{\mathbb{F}}^2/n \\ & = O_P(n^{-\frac{1}{2}}) \|\mathbf{F}\|_{\mathbb{F}}^2/n. \end{aligned}$$

Using $\|(\mathbf{I}_T - \widehat{\mathbf{A}})^{-1}\|_{\text{op}} = O_P(1)$ again, it follows

$$\|\widehat{\mathbf{S}} - \mathbf{S}\|_{\mathbb{F}} \leq O_P(n^{-\frac{1}{2}}) \|\mathbf{F}\|_{\mathbb{F}}^2/n. \quad (33)$$

A similar argument leads to

$$\|\widehat{\mathbf{S}} - \mathbf{S}\|_{\mathbb{F}} \leq O_P(n^{-\frac{1}{2}}) (\text{Tr}(\mathbf{S}) + \|\mathbf{H}\|_{\mathbb{F}}^2). \quad (34)$$

■

Proof of Corollary 3.6. Under Assumption 4(i) and 5, we proceed to bound $\|\mathbf{F}\|_{\mathbb{F}}^2/n$ in terms of $\text{Tr}(\mathbf{S})$. Let $L(\mathbf{B}) = \frac{1}{2n} \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_{\mathbb{F}}^2 + \lambda \|\mathbf{B}\|_{2,1} + \frac{\tau}{2} \|\mathbf{B}\|_{\mathbb{F}}^2$ be the objective function in (6), then $L(\widehat{\mathbf{B}}) \leq L(\mathbf{0})$ by definition of $\widehat{\mathbf{B}}$. Thus,

$$\frac{1}{2n} \|\mathbf{F}\|_{\mathbb{F}}^2 \leq \frac{1}{2n} \|\mathbf{F}\|_{\mathbb{F}}^2 + \lambda \|\widehat{\mathbf{B}}\|_{2,1} + \frac{\tau}{2} \|\widehat{\mathbf{B}}\|_{\mathbb{F}}^2 \leq \frac{1}{2n} \|\mathbf{Y}\|_{\mathbb{F}}^2.$$

Now we bound $\frac{1}{n} \|\mathbf{Y}\|_{\mathbb{F}}^2$ by Hanson-Wright inequality. Since $\mathbf{Y} = \mathbf{X}\mathbf{B}^* + \mathbf{E}$, the rows of \mathbf{Y} are i.i.d. $\mathcal{N}_T(\mathbf{0}, \Sigma_{\mathbf{y}})$ with $\Sigma_{\mathbf{y}} = (\mathbf{B}^*)^\top \Sigma \mathbf{B}^* + \mathbf{S}$, then $\text{vec}(\mathbf{Y}^\top) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n \otimes \Sigma_{\mathbf{y}})$, and $\xi \stackrel{\text{def}}{=} [\mathbf{I}_n \otimes \Sigma_{\mathbf{y}}]^{-\frac{1}{2}} \text{vec}(\mathbf{Y}^\top) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{nT})$. Since $\|\mathbf{Y}\|_{\mathbb{F}}^2 = [\text{vec}(\mathbf{Y}^\top)]^\top \text{vec}(\mathbf{Y}^\top) = \xi^\top (\mathbf{I}_n \otimes \Sigma_{\mathbf{y}}) \xi$, we apply the following variant of Hanson-Wright inequality.

Lemma B.6 (Lemma 1 in Laurent and Massart (2000)). For $\xi \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_N)$ and matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$, the following inequality holds for any $x > 0$,

$$\mathbb{P}(\xi^\top \mathbf{A} \xi - \text{Tr}(\mathbf{A}) \leq 2\sqrt{x} \|\mathbf{A}\|_{\mathbb{F}} + 2x \|\mathbf{A}\|_{\text{op}}) \geq 1 - \exp(-x).$$

In our case, we take $\mathbf{A} = (\mathbf{I}_n \otimes \Sigma_{\mathbf{y}})$, then $\text{Tr}(\mathbf{A}) = n \text{Tr}(\Sigma_{\mathbf{y}})$, $\|\mathbf{A}\|_{\mathbb{F}} = \sqrt{n} \|\Sigma_{\mathbf{y}}\|_{\mathbb{F}} \leq \sqrt{n} \text{Tr}(\Sigma_{\mathbf{y}})$, $\|\mathbf{A}\|_{\text{op}} = \|\Sigma_{\mathbf{y}}\|_{\text{op}} \leq \text{Tr}(\Sigma_{\mathbf{y}})$, thus with probability at least $1 - \exp(-x)$,

$$\|\mathbf{Y}\|_{\mathbb{F}}^2 - n \text{Tr}(\Sigma_{\mathbf{y}}) \leq 2\sqrt{nx} \text{Tr}(\Sigma_{\mathbf{y}}) + 2x \text{Tr}(\Sigma_{\mathbf{y}}).$$

Take $x = n$, then with probability at least $1 - \exp(-n)$,

$$\|\mathbf{F}\|_{\mathbb{F}}^2/n \leq \|\mathbf{Y}\|_{\mathbb{F}}^2/n \leq 5 \text{Tr}(\Sigma_{\mathbf{y}}).$$

Thus, $\|\mathbf{F}\|_{\text{F}}^2/n = O_P(1) \text{Tr}(\boldsymbol{\Sigma}_{\mathbf{y}})$. Together with (33), we obtain

$$\|\widehat{\mathbf{S}} - \mathbf{S}\|_{\text{F}} \leq O_P(n^{-\frac{1}{2}}) \text{Tr}(\boldsymbol{\Sigma}_{\mathbf{y}}).$$

Note that by Assumption 5, $\text{Tr}(\boldsymbol{\Sigma}_{\mathbf{y}}) = \|\boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{B}^*\|_{\text{F}}^2 + \text{Tr}(\mathbf{S}) \leq (1 + \varsigma \text{nr}) \text{Tr}(\mathbf{S})$. Therefore, we obtain

$$\|\widehat{\mathbf{S}} - \mathbf{S}\|_{\text{F}} \leq O_P(n^{-\frac{1}{2}}) \text{Tr}(\mathbf{S}).$$

Furthermore, since $\text{Tr}(\mathbf{S}) \leq \sqrt{T} \|\mathbf{S}\|_{\text{F}}$ and $T = o(n)$, we have

$$\|\widehat{\mathbf{S}} - \mathbf{S}\|_{\text{F}} \leq O_P(\sqrt{T/n}) \|\mathbf{S}\|_{\text{F}} = o_P(1) \|\mathbf{S}\|_{\text{F}}.$$

Finally, since $\|\mathbf{S}\|_* = \text{Tr}(\mathbf{S})$, by the triangular inequality

$$|\|\widehat{\mathbf{S}}\|_* - \text{Tr}(\mathbf{S})| \leq \|\widehat{\mathbf{S}} - \mathbf{S}\|_* \leq \sqrt{T} \|\widehat{\mathbf{S}} - \mathbf{S}\|_{\text{F}} \leq O_P(\sqrt{T/n}) \text{Tr}(\mathbf{S}) = o_P(1) \text{Tr}(\mathbf{S}).$$

■

Proof of Corollary 3.7. For $\tau = 0$, by the optimality of $\widehat{\mathbf{B}}$ in (6),

$$\frac{1}{2n} \|\mathbf{F}\|_{\text{F}}^2 + \lambda \|\widehat{\mathbf{B}}\|_{2,1} \leq \frac{1}{2n} \|\mathbf{E}\|_{\text{F}}^2 + \lambda \|\mathbf{B}^*\|_{2,1}.$$

Note that $\mathbf{F} = \mathbf{E} - \mathbf{X}(\widehat{\mathbf{B}} - \mathbf{B}^*) = \mathbf{E} - \mathbf{Z}\mathbf{H}$, expanding the squares and rearranging terms yields

$$\|\mathbf{Z}\mathbf{H}\|_{\text{F}}^2 \leq 2\langle \mathbf{E}, \mathbf{Z}\mathbf{H} \rangle + 2n\lambda(\|\mathbf{B}^*\|_{2,1} - \|\widehat{\mathbf{B}}\|_{2,1}) \leq 2\langle \mathbf{E}, \mathbf{Z}\mathbf{H} \rangle + 2n\lambda \|\widehat{\mathbf{B}} - \mathbf{B}^*\|_{2,1}. \quad (35)$$

From assumptions in this corollary, $\widehat{\mathbf{B}} - \mathbf{B}^*$ has at most $(1 - c)n$ rows. Thus, in the event U_2 , we have

$$n\eta \|\mathbf{H}\|_{\text{F}}^2 = n\eta \|\boldsymbol{\Sigma}^{1/2}(\widehat{\mathbf{B}} - \mathbf{B}^*)\|_{\text{F}}^2 \leq \|\mathbf{X}(\widehat{\mathbf{B}} - \mathbf{B}^*)\|_{\text{F}}^2 = \|\mathbf{Z}\mathbf{H}\|_{\text{F}}^2.$$

We bound the right-hand side two terms in (35) by Cauchy-Schwarz inequality,

$$\|\widehat{\mathbf{B}} - \mathbf{B}^*\|_{2,1} \leq \sqrt{(1 - c)n} \|\widehat{\mathbf{B}} - \mathbf{B}^*\|_{\text{F}} \leq \frac{\sqrt{(1 - c)n}}{\sqrt{\phi_{\min}(\boldsymbol{\Sigma})}} \|\mathbf{H}\|_{\text{F}} \leq \frac{\sqrt{1 - c}}{\sqrt{\eta \phi_{\min}(\boldsymbol{\Sigma})}} \|\mathbf{Z}\mathbf{H}\|_{\text{F}},$$

and $\langle \mathbf{E}, \mathbf{Z}\mathbf{H} \rangle \leq \|\mathbf{E}\|_{\text{F}} \|\mathbf{Z}\mathbf{H}\|_{\text{F}} \leq \|\mathbf{S}^{\frac{1}{2}}\|_{\text{F}} \|\mathbf{E}\mathbf{S}^{-\frac{1}{2}}\|_{\text{op}} \|\mathbf{Z}\mathbf{H}\|_{\text{F}}$.

Therefore, by canceling a factor $\|\mathbf{Z}\mathbf{H}\|_{\text{F}}$ from both sides of (35), we have

$$\sqrt{n\eta} \|\mathbf{H}\|_{\text{F}} \leq \|\mathbf{Z}\mathbf{H}\|_{\text{F}} \leq 2\|\mathbf{S}^{\frac{1}{2}}\|_{\text{F}} \|\mathbf{E}\mathbf{S}^{-\frac{1}{2}}\|_{\text{op}} + \frac{2\sqrt{(1 - c)n}\lambda}{\sqrt{\eta \phi_{\min}(\boldsymbol{\Sigma})}}.$$

Using $(a + b)^2 \leq 2a^2 + 2b^2$,

$$\|\mathbf{H}\|_{\text{F}}^2 \leq \frac{4}{n\eta} \text{Tr}(\mathbf{S}) \|\mathbf{E}\mathbf{S}^{-\frac{1}{2}}\|_{\text{op}}^2 + \frac{4(1 - c)n\lambda^2}{\eta^2 \phi_{\min}(\boldsymbol{\Sigma})}.$$

Hence, using λ is of the form $\mu\sqrt{\text{Tr}(\mathbf{S})/n}$, we have

$$\begin{aligned} & \text{Tr}(\mathbf{S}) + \|\mathbf{H}\|_{\text{F}}^2 \\ & \leq (1 + 4\eta^{-1}n^{-1}\|\mathbf{E}\mathbf{S}^{-\frac{1}{2}}\|_{\text{op}}^2) \text{Tr}(\mathbf{S}) + \frac{4(1-c)\mu^2}{\eta^2\phi_{\min}(\boldsymbol{\Sigma})} \text{Tr}(\mathbf{S}) \\ & \leq O_P(1)(1 + \mu^2) \text{Tr}(\mathbf{S}), \end{aligned}$$

where we used that $n^{-1}\|\mathbf{E}\mathbf{S}^{-\frac{1}{2}}\|_{\text{op}} = O_P(1)$ by (Davidson and Szarek, 2001, Theorem II.13) and $T = o(n)$. Now, by Theorem 3.5,

$$\|\widehat{\mathbf{S}} - \mathbf{S}\|_{\text{F}} \leq O_P(n^{-\frac{1}{2}})[\text{Tr}(\mathbf{S}) + \|\mathbf{H}\|_{\text{F}}^2] \leq O_P(n^{-\frac{1}{2}})(1 + \mu^2) \text{Tr}(\mathbf{S}),$$

where the $O_P(\cdot)$ hides constants depending on $\gamma, c, \phi_{\min}(\boldsymbol{\Sigma})$ since η is a constant that only depends on γ, c . ■

Proof of Theorem 3.10. From the definitions of $\mathbf{Q}_2, \mathbf{Q}_3$ in Propositions B.4 and B.5, we have

$$\begin{aligned} & \mathbf{Q}_2 + \mathbf{Q}_2^{\top} + \mathbf{Q}_3 \\ & = \frac{n^{-2}(\mathbf{F}^{\top}\mathbf{Z}\mathbf{Z}^{\top}\mathbf{F} + \widehat{\mathbf{A}}\mathbf{F}^{\top}\mathbf{F} + \mathbf{F}^{\top}\mathbf{F}\widehat{\mathbf{A}} - p\mathbf{F}^{\top}\mathbf{F} - (n\mathbf{I}_T - \widehat{\mathbf{A}})\mathbf{H}^{\top}\mathbf{H}(n\mathbf{I}_T - \widehat{\mathbf{A}})}{(\|\mathbf{H}\|_{\text{F}}^2 + \|\mathbf{F}\|_{\text{F}}^2/n)n^{-\frac{1}{2}}}. \end{aligned}$$

Therefore,

$$\begin{aligned} & \|(I_T - \widehat{\mathbf{A}}/n)\mathbf{H}^{\top}\mathbf{H}(I_T - \widehat{\mathbf{A}}/n) - n^{-2}(\mathbf{F}^{\top}\mathbf{Z}\mathbf{Z}^{\top}\mathbf{F} + \widehat{\mathbf{A}}\mathbf{F}^{\top}\mathbf{F} + \mathbf{F}^{\top}\mathbf{F}\widehat{\mathbf{A}} - p\mathbf{F}^{\top}\mathbf{F})\|_{\text{F}} \\ & = \|\mathbf{Q}_2 + \mathbf{Q}_2^{\top} + \mathbf{Q}_3\|_{\text{F}}(\|\mathbf{H}\|_{\text{F}}^2 + \|\mathbf{F}\|_{\text{F}}^2/n)n^{-\frac{1}{2}} \\ & \leq \Theta_3(\|\mathbf{H}\|_{\text{F}}^2 + \|\mathbf{F}\|_{\text{F}}^2/n)n^{-\frac{1}{2}}, \end{aligned}$$

where $\Theta_3 = 2\|\mathbf{Q}_2\|_{\text{F}} + \|\mathbf{Q}_3\|_{\text{F}}$. The conclusion thus follows by Propositions B.4 and B.5. ■

Acknowledgments

The authors would like to thank the anonymous referees, an Associate Editor and the Editor for their constructive comments that improved the quality of this paper. The authors acknowledge the Office of Advanced Research Computing (OARC) at Rutgers, The State University of New Jersey for providing access to the Amarel cluster and associated research computing resources that have contributed to the results reported here. URL: <https://oarc.rutgers.edu>.

Funding

Pierre C. Bellec is partially supported by the NSF Grants DMS-1811976 and DMS-1945428. Gabriel Romon is supported by a PhD scholarship from CREST.

Supplementary Material

Supplement I

Additional simulation results and proofs of results in Appendix B.1.

Supplement II

Code for reproducing the numerical experiments in the main paper and the Supplement I.

References

- BAYATI, M., ERDOĞDU, M. A. and MONTANARI, A. (2013). Estimating lasso risk and noise level. *Adv. Neural Inf. Process. Syst.* **26**.
- BAYATI, M. and MONTANARI, A. (2012). The LASSO risk for Gaussian matrices. *IEEE Trans. Inform. Theory* **58** 1997–2017.
- BELLEÇ, P. C. (2020). Out-of-sample error estimate for robust M-estimators with convex penalty. *arXiv preprint arXiv:2008.11840*.
- BELLEÇ, P. and KUCHIBHOTLA, A. (2019). First order expansion of convex regularized estimators. *Adv. Neural Inf. Process. Syst.* **32**.
- BELLEÇ, P. C. and ROMON, G. (2021). Chi-square and normal inference in high-dimensional multi-task regression. *arXiv preprint arXiv:2107.07828*.
- BELLEÇ, P. and TSYBAKOV, A. (2016). Bounds on the prediction error of penalized least squares estimators with convex penalty. In *International Conference on Modern Problems of Stochastic Analysis and Statistics* 315–333. Springer.
- BELLEÇ, P. C. and ZHANG, C.-H. (2019). De-biasing convex regularized estimators and interval estimation in linear models. *arXiv preprint arXiv:1912.11943*.
- BELLEÇ, P. C. and ZHANG, C.-H. (2021). Second-order Stein: SURE for SURE and other applications in high-dimensional inference. *Ann. Statist.* **49** 1864–1903.
- BELLONI, A., CHERNOZHUKOV, V. and WANG, L. (2011). Square-root lasso: pivotal recovery of sparse signals via conic programming. *Biometrika* **98** 791–806.
- BELLONI, A., CHERNOZHUKOV, V. and WANG, L. (2014). Pivotal estimation via square-root Lasso in nonparametric regression. *Ann. Statist.* **42** 757–788.
- BERTRAND, Q., MASSIAS, M., GRAMFORT, A. and SALMON, J. (2019). Handling correlated and repeated measurements with the smoothed multivariate square-root Lasso. *Adv. Neural Inf. Process. Syst.* **32**.
- BICKEL, P. J. and LEVINA, E. (2008). Covariance regularization by thresholding. *Ann. Statist.* **36** 2577–2604.
- BOUCHERON, S., LUGOSI, G. and MASSART, P. (2013). *Concentration inequalities: A nonasymptotic theory of independence*. Oxford University Press.
- CAI, T. T., ZHANG, C.-H. and ZHOU, H. H. (2010). Optimal rates of convergence for covariance matrix estimation. *Ann. Statist.* **38** 2118–2144.
- CELENTANO, M. and MONTANARI, A. (2021). CAD: Debiasing the Lasso with inaccurate covariate model. *arXiv preprint arXiv:2107.14172*.
- CHEN, S. and BANERJEE, A. (2017). Alternating estimation for structured high-dimensional multi-response models. *Adv. Neural Inf. Process. Syst.* **30**.
- DAVIDSON, K. R. and SZAREK, S. J. (2001). Local operator theory, random matrices and Banach spaces. *Handbook of the geometry of Banach spaces* **1** 131.
- DICKER, L. H. (2014). Variance estimation in high-dimensional linear models. *Biometrika* **101** 269–284.
- DICKER, L. H. and ERDOĞDU, M. A. (2016). Maximum likelihood for variance estimation in high-dimensional linear models. In *Artificial Intelligence and Statistics* 159–167. PMLR.
- DOBRIAN, E. and WAGER, S. (2018). High-dimensional asymptotics of prediction: ridge regression and classification. *Ann. Statist.* **46** 247–279.
- EL KAROUI, N. (2008). Operator norm consistent estimation of large-dimensional sparse covariance matrices. *Ann. Statist.* **36** 2717–2756.

- FAN, J., GUO, S. and HAO, N. (2012). Variance estimation using refitted cross-validation in ultrahigh dimensional regression. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **74** 37–65.
- FOURDRINIER, D., HADDOUCHE, A. M. and MEZOUE, F. (2021). Covariance matrix estimation under data-based loss. *Statist. Probab. Lett.* **177** Paper No. 109160, 7.
- FRIEDMAN, J. H., HASTIE, T. and TIBSHIRANI, R. (2010). Regularization Paths for Generalized Linear Models via Coordinate Descent. *J. Stat. Softw.* **33** 1–22.
- GEER, S. v. D. and STUCKY, B. (2016). χ^2 -confidence sets in high-dimensional regression. In *Statistical analysis for high-dimensional data* 279–306. Springer.
- JANSON, L., FOYGE BARBER, R. and CANDÈS, E. (2017). EigenPrism: inference for high dimensional signal-to-noise ratios. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **79** 1037–1065.
- KOLTCHINSKII, V. and LOUNICI, K. (2017). Concentration inequalities and moment bounds for sample covariance operators. *Bernoulli* **23** 110–133.
- LAURENT, B. and MASSART, P. (2000). Adaptive estimation of a quadratic functional by model selection. *Ann. Statist.* **28** 1302–1338.
- LECUÉ, G. and MENDELSON, S. (2017). Sparse recovery under weak moment assumptions. *J. Eur. Math. Soc. (JEMS)* **19** 881–904.
- LIU, H. and ZHANG, J. (2009). Estimation consistency of the group lasso and its applications. In *Artificial Intelligence and Statistics* 376–383. PMLR.
- LOUNICI, K., PONTIL, M., VAN DE GEER, S. and TSYBAKOV, A. B. (2011). Oracle inequalities and optimal inference under group sparsity. *Ann. Statist.* **39** 2164–2204.
- MIOLANE, L. and MONTANARI, A. (2021). The distribution of the Lasso: uniform control over sparse balls and adaptive parameter tuning. *Ann. Statist.* **49** 2313–2335.
- MOLSTAD, A. J. (2022). New Insights for the Multivariate Square-Root Lasso. *J. Mach. Learn. Res.* **23** 1–52.
- OBOZINSKI, G., WAINWRIGHT, M. J. and JORDAN, M. I. (2011). Support union recovery in high-dimensional multivariate regression. *Ann. Statist.* **39** 1–47.
- PEDREGOSA, F., VAROQUAUX, G., GRAMFORT, A., MICHEL, V., THIRION, B., GRISEL, O., BLONDEL, M., PRETTENHOFER, P., WEISS, R., DUBOURG, V., VANDERPLAS, J., PASSOS, A., COURNAPEAU, D., BRUCHER, M., PERROT, M. and DUCHESNAY, E. (2011). Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* **12** 2825–2830.
- SIMON, N., FRIEDMAN, J. and HASTIE, T. (2013). A blockwise descent algorithm for group-penalized multiresponse and multinomial regression. *arXiv preprint arXiv:1311.6529*.
- STEIN, C. M. (1981). Estimation of the mean of a multivariate normal distribution. *Ann. Statist.* 1135–1151.
- SUN, T. and ZHANG, C.-H. (2012). Scaled sparse linear regression. *Biometrika* **99** 879–898.
- TAN, K., ROMON, G. and BELLEC, P. C. (2023). Supplement to "Noise covariance estimation in multi-task high-dimensional linear models".
- YU, G. and BIEN, J. (2019). Estimating the error variance in a high-dimensional linear model. *Biometrika* **106** 533–546.

Noise Covariance Estimation in Multi-Task High-dimensional Linear Models Supplementary Material

Kai Tan, Gabriel Romon, Pierre C. Bellec

In this supplement, we first present additional simulation results for estimating two types of noise covariance with different (n, p) in Section C. To complete the proofs of all theoretical results in the paper, it remains to prove Lemma B.1, B.2 and Propositions B.3 to B.5 in Appendix B of the paper. We prove Lemmas B.1 and B.2 in Section D, and present in Section E a few useful lemmas before proving Propositions B.3 to B.5 in Section F. The proofs of lemmas in Section E are deferred to Section G.

C. Additional simulation results

While Figures 2 and 3 in the main paper present heatmaps for estimating full and low-rank noise covariance S with $n = 1000$, this section provides additional estimating results from $n = 1500$ and $n = 2000$ (both with $p = 1.5n$).

C.1. Heatmaps for estimating full rank S with $n = 1500, 2000$

When estimating the full rank noise covariance matrix S with (t, t') -th entry $S_{t,t'} = \frac{\cos(t-t')}{1+\sqrt{|t-t'|}}$, the heatmaps for different estimators from $n = 1500$ and $n = 2000$ are presented in Figures 4 and 5, respectively. The comparisons between four methods in Figures 4 and 5 are similar to the case $n = 1000$ in Figure 2 of the main paper; our proposed estimator outperforms the naive estimator and method-of-moments estimator.

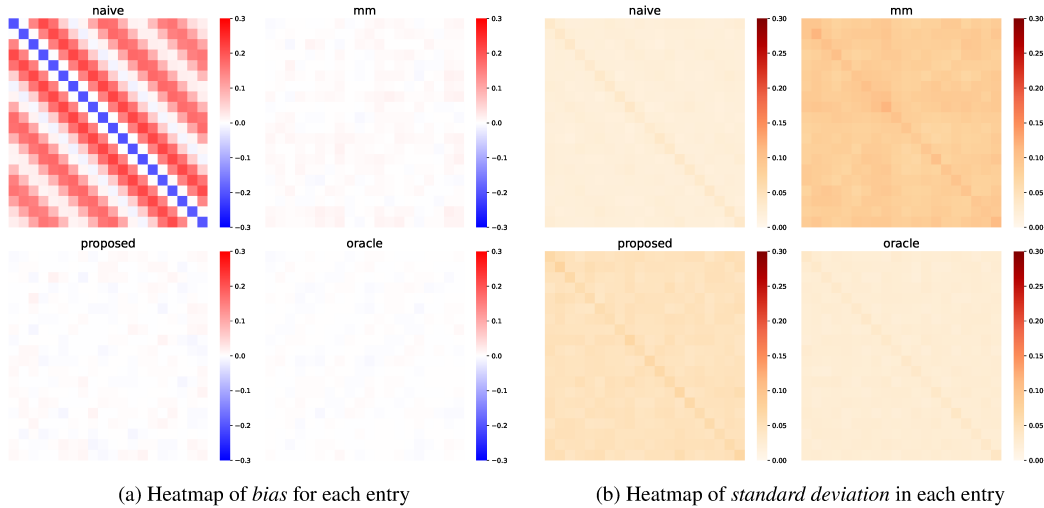


Figure 4: Heatmaps for estimation of full rank S with $n = 1500$ over 100 repetitions.

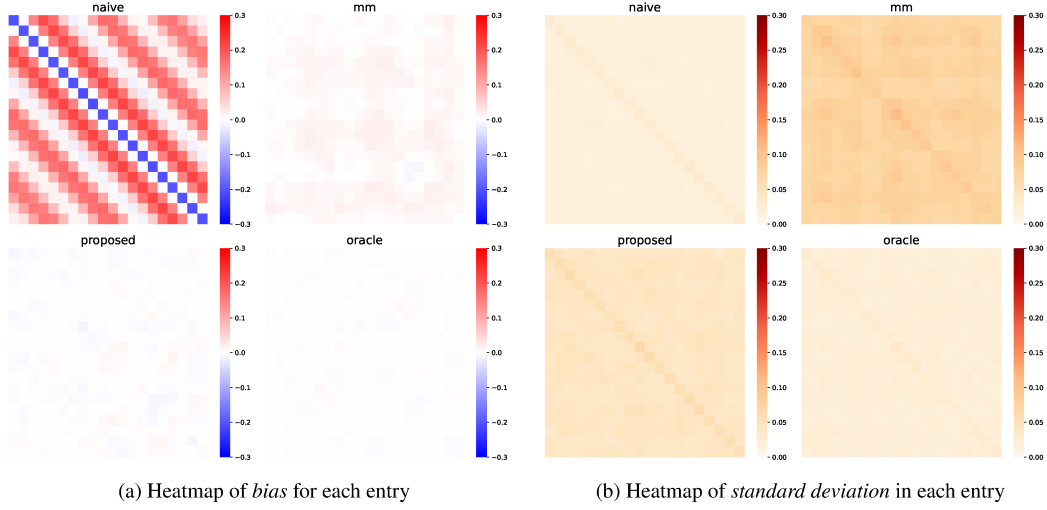


Figure 5: Heatmaps for estimation of full rank S with $n = 2000$ over 100 repetitions.

C.2. Heatmaps for estimating low rank S with $n = 1500, 2000$

When estimating the low rank noise covariance matrix with $S = uu^\top$, and $u \in \mathbb{R}^{T \times 10}$ with entries are i.i.d. from $\mathcal{N}(0, 1/T)$. We present the heatmaps for different estimators with 1500, 2000 in Figures 6 and 7 below. The comparisons between four methods in Figures 6 and 7 are similar to the case $n = 1000$ in Figure 3 of the main paper. All of these figures convince us that besides the oracle estimator, the proposed estimator has the best performance.

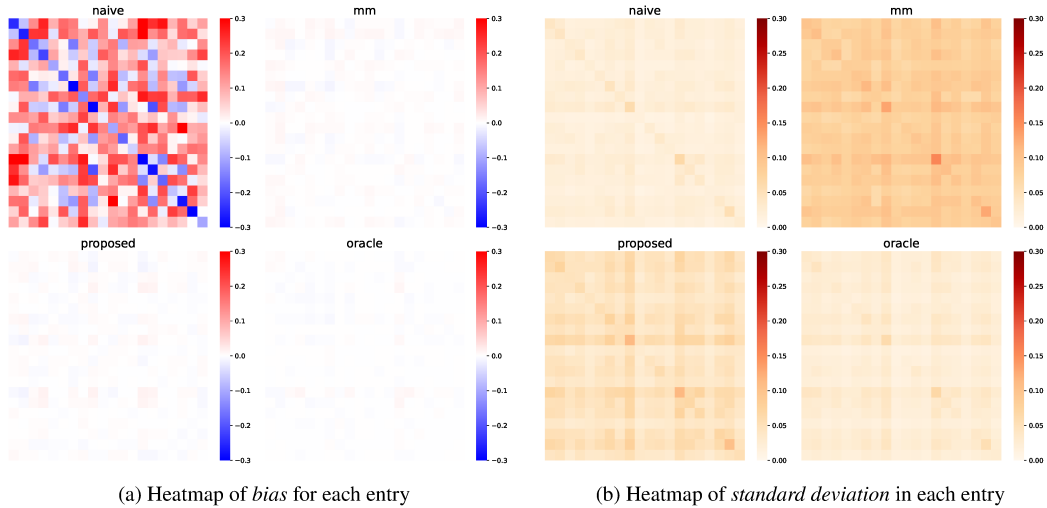
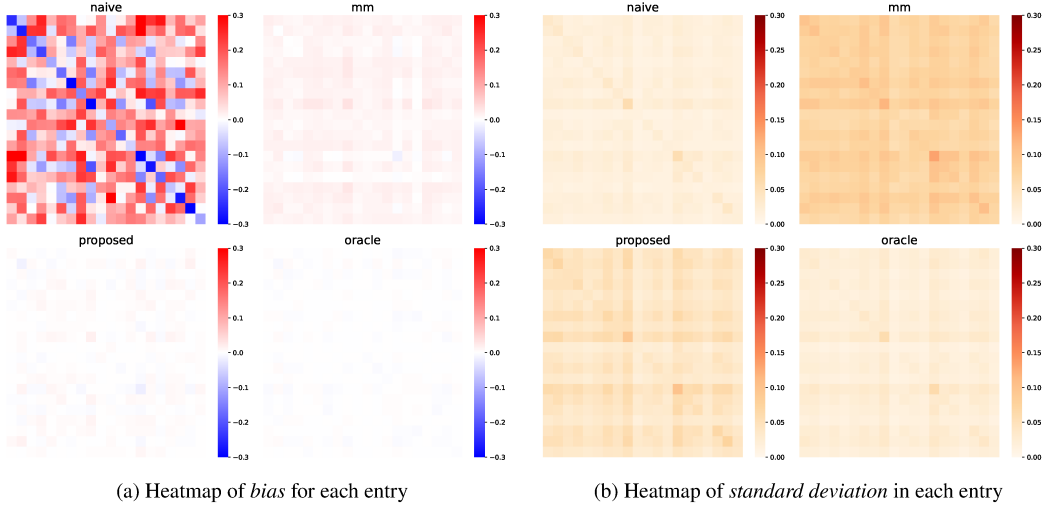


Figure 6: Heatmaps for estimation of low rank S with $n = 1500$ over 100 repetitions.

Figure 7: Heatmaps for estimation of low rank \mathbf{S} with $n = 2000$ over 100 repetitions.

D. Proofs of Lemmas B.1 and B.2

Proof of Lemma B.1. (i) For any $\mathbf{u} \in \mathbb{R}^T$, by definition (10),

$$\begin{aligned}
 \mathbf{u}^\top \widehat{\mathbf{A}} \mathbf{u} &= \text{Tr}[(\mathbf{u}^\top \otimes \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger (\mathbf{u} \otimes \mathbf{X}_{\hat{\mathcal{J}}}^\top)] \\
 &\leq \text{Tr}[(\mathbf{u}^\top \otimes \mathbf{X}_{\hat{\mathcal{J}}}) [\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}})^\dagger] (\mathbf{u} \otimes \mathbf{X}_{\hat{\mathcal{J}}}^\top)] \\
 &= \text{Tr}[(\mathbf{u}^\top \mathbf{I}_T \mathbf{u}) \otimes [\mathbf{X}_{\hat{\mathcal{J}}} (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}})^\dagger \mathbf{X}_{\hat{\mathcal{J}}}^\top]] \\
 &= \|\mathbf{u}\|^2 \text{Tr}[\mathbf{X}_{\hat{\mathcal{J}}} (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}})^\dagger \mathbf{X}_{\hat{\mathcal{J}}}^\top] \\
 &= \|\mathbf{u}\|^2 \text{Tr}[\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}})^\dagger].
 \end{aligned}$$

Let $r = \text{rank}(\mathbf{X}_{\hat{\mathcal{J}}}) \leq \min(n, |\hat{\mathcal{J}}|)$ be the rank of $\mathbf{X}_{\hat{\mathcal{J}}}$, and $\widehat{\phi}_1 \geq \dots \geq \widehat{\phi}_r > 0$ be the nonzero eigenvalues of $\frac{1}{n} \mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}}$. We have

$$\begin{aligned}
 \|\widehat{\mathbf{A}}/n\|_{\text{op}} &\leq \frac{1}{n} \text{Tr}[\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}})^\dagger] \\
 &= \frac{1}{n} \text{Tr}[\frac{1}{n} \mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} (\frac{1}{n} \mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + \tau \mathbf{P}_{\hat{\mathcal{J}}})^\dagger] \\
 &\leq \frac{r}{n} \left\| \frac{1}{n} \mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} (\frac{1}{n} \mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + \tau \mathbf{P}_{\hat{\mathcal{J}}})^\dagger \right\|_{\text{op}} \\
 &\leq \frac{\widehat{\phi}_1}{\widehat{\phi}_1 + \tau} \leq 1.
 \end{aligned}$$

Thus, $\|\mathbf{I} - \widehat{\mathbf{A}}/n\|_{\text{op}} \leq 1$ as $\widehat{\mathbf{A}}$ is positive semi-definite.

(ii) Note that $\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1}\|_{\text{op}} = (1 - \|\widehat{\mathbf{A}}/n\|_{\text{op}})^{-1} \leq 1 + \frac{\widehat{\phi}_1}{\tau}$, where $\widehat{\phi}_1 = \|\frac{1}{n}\mathbf{X}^\top \mathbf{X}_{\widehat{\mathcal{J}}}\|_{\text{op}} \leq \|\frac{1}{n}\mathbf{X}^\top \mathbf{X}\|_{\text{op}} \leq \frac{1}{n}\|\mathbf{X}^\top \mathbf{\Sigma}^{-\frac{1}{2}}\|_{\text{op}}^2 \|\mathbf{\Sigma}\|_{\text{op}}$. Therefore,

(1) in the event $\{\|\mathbf{X}\mathbf{\Sigma}^{-\frac{1}{2}}\|_{\text{op}} < 2\sqrt{n} + \sqrt{p}\}$, we have

$$\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1}\|_{\text{op}} \leq 1 + \tau^{-1}(2 + \sqrt{p/n})^2 \|\mathbf{\Sigma}\|_{\text{op}} = 1 + (\tau')^{-1}(2 + \sqrt{p/n})^2.$$

(2) $\mathbb{E}[\widehat{\phi}_1] \leq \mathbb{E}[n^{-1}\|\mathbf{X}^\top \mathbf{\Sigma}^{-\frac{1}{2}}\|_{\text{op}}^2 \|\mathbf{\Sigma}\|_{\text{op}}] \leq [(1 + \sqrt{p/n})^2 + n^{-1}]\|\mathbf{\Sigma}\|_{\text{op}}$ by (36). Hence,

$$\mathbb{E}\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1}\|_{\text{op}} \leq 1 + \tau^{-1}\mathbb{E}[\widehat{\phi}_1] \leq 1 + (\tau')^{-1}[(1 + \sqrt{p/n})^2 + n^{-1}].$$

■

Proof of Lemma B.2. (i) For $\tau = 0$, using the same argument as proof of Lemma B.1, we obtain

$$\mathbf{u}^\top \widehat{\mathbf{A}} \mathbf{u} \leq \|\mathbf{u}\|^2 \text{Tr}[\mathbf{X}_{\widehat{\mathcal{J}}}^\top \mathbf{X}_{\widehat{\mathcal{J}}} (\mathbf{X}_{\widehat{\mathcal{J}}}^\top \mathbf{X}_{\widehat{\mathcal{J}}})^\dagger] \leq \|\mathbf{u}\|^2 |\widehat{\mathcal{J}}|.$$

Thus, in the event U_1 , we have $\|\widehat{\mathbf{A}}\|_{\text{op}}/n \leq |\widehat{\mathcal{J}}|/n \leq (1 - c)/2 < 1$, hence

$$\|\mathbf{I}_T - \widehat{\mathbf{A}}/n\|_{\text{op}} \leq 1.$$

(ii) In the event U_1 , we have $\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1}\|_{\text{op}} = (1 - \|\widehat{\mathbf{A}}/n\|_{\text{op}})^{-1} \leq (1 - (1 - c)/2)^{-1}$. Furthermore, $\mathbb{E}[I(U_1)\|(\mathbf{I}_T - \widehat{\mathbf{A}}/n)^{-1}\|_{\text{op}}] \leq (1 - (1 - c)/2)^{-1}$.

■

E. Preliminary results needed for proving Propositions B.3 to B.5

Let us first introduce two events besides the event $U_1 = \{\|\widehat{\mathbf{B}}\|_0 \leq n(1 - c)/2\}$ in Assumption 4(ii), we define events U_2 and U_3 as below,

$$U_2 = \left\{ \inf_{\mathbf{b} \in \mathbb{R}^p: \|\mathbf{b}\|_0 \leq (1-c)n} \|\mathbf{X}\mathbf{b}\|^2 / (n\|\mathbf{\Sigma}^{\frac{1}{2}}\mathbf{b}\|^2) > \eta \right\},$$

$$U_3 = \left\{ \|\mathbf{X}\mathbf{\Sigma}^{-\frac{1}{2}}\|_{\text{op}} < 2\sqrt{n} + \sqrt{p} \right\}.$$

Under Assumptions 2 and 3, (Bellec and Zhang, 2019, Lemma B.1) guarantees $\mathbb{P}(U_2) \geq 1 - C(\gamma, c)e^{-C(\gamma, c)n}$ for some constant η that only depends on constants γ, c . Under Assumption 2, (Davidson and Szarek, 2001, Theorem II.13) guarantees $\mathbb{P}(U_3) > 1 - e^{-n/2}$ and there exists a random variable $z \sim \mathcal{N}(0, 1)$ s.t. $\|\mathbf{X}\mathbf{\Sigma}^{-\frac{1}{2}}\|_{\text{op}} \leq \sqrt{n} + \sqrt{p} + z$ almost surely. Therefore, under Assumptions 2 and 3, we have

$$\mathbb{E}[\|n^{-\frac{1}{2}}\mathbf{X}\mathbf{\Sigma}^{-\frac{1}{2}}\|_{\text{op}}^2] \leq (1 + \sqrt{p/n})^2 + n^{-1} \leq C(\gamma). \quad (36)$$

Furthermore, under Assumptions 2, 3 and 4(ii), $\mathbb{P}(U_1 \cap U_2 \cap U_3) \rightarrow 1$ by a union bound, and for large enough n ,

$$\begin{aligned} \mathbb{P}\{(U_1 \cap U_2 \cap U_3)^c\} &< 1/T + C(\gamma, c)e^{-n/C(\gamma, c)} \\ &= \frac{1}{T}(1 + TC(\gamma, c)e^{-n/C(\gamma, c)}) \end{aligned}$$

$$\begin{aligned}
&< \frac{1}{T} (1 + C(\gamma, c) e^{\sqrt{n}-n/C(\gamma, c)}) \\
&< \frac{1}{T} C(\gamma, c).
\end{aligned} \tag{37}$$

E.1. Lipschitz and differential properties for a given, fixed noise matrix E

We need to study Lipschitz and differential properties of certain mappings when the noise matrix E is fixed. Let $g : \mathbb{R}^{p \times T} \rightarrow \mathbb{R}$ defined by $g(\mathbf{B}) = \tau \|\mathbf{B}\|_F^2/2 + \lambda \|\mathbf{B}\|_{2,1}$ be the penalty in (6). For a fixed value of E , define the mappings

$$\mathbf{Z} \mapsto \mathbf{H}(\mathbf{Z}) = \arg \min_{\mathbf{H} \in \mathbb{R}^{p \times T}} \frac{1}{2n} \|\mathbf{E} - \mathbf{Z}\mathbf{H}\|_F^2 + g(\mathbf{\Sigma}^{-1/2} \mathbf{H}) \quad (\mathbb{R}^{n \times p} \rightarrow \mathbb{R}^{p \times T}) \tag{38}$$

$$\mathbf{Z} \mapsto \mathbf{F}(\mathbf{Z}) = \mathbf{E} - \mathbf{Z}\mathbf{H}(\mathbf{Z}) \quad (\mathbb{R}^{n \times p} \rightarrow \mathbb{R}^{n \times T}) \tag{39}$$

$$\mathbf{Z} \mapsto D(\mathbf{Z}) = (\|\mathbf{H}(\mathbf{Z})\|_F^2 + \|\mathbf{F}(\mathbf{Z})\|_F^2/n)^{1/2} \quad (\mathbb{R}^{n \times p} \rightarrow \mathbb{R}). \tag{40}$$

Next, define the random variable $\mathbf{Z} = \mathbf{X}\mathbf{\Sigma}^{-\frac{1}{2}} \in \mathbb{R}^{n \times p}$, and let us use the convention that if arguments of \mathbf{H}, \mathbf{F} or D are omitted then these mappings are implicitly taken at the realized value of the random variable $\mathbf{Z} = \mathbf{X}\mathbf{\Sigma}^{-\frac{1}{2}} \in \mathbb{R}^{n \times p}$ where \mathbf{X} is the observed design matrix. With this convention and by definition of the above mappings, we then have $\mathbf{H} = \mathbf{H}(\mathbf{Z}) = \mathbf{\Sigma}^{1/2}(\widehat{\mathbf{B}} - \mathbf{B}^*)$ as well as $\mathbf{F} = \mathbf{F}(\mathbf{Z}) = \mathbf{Y} - \mathbf{X}\widehat{\mathbf{B}}$ and $D = [\|\mathbf{H}\|_F^2 + \|\mathbf{F}\|_F^2/n]^{1/2}$ so that the notation is consistent with the rest of the paper (in particular, with (9) in the main paper).

Finally, denote the (i, j) -th entry of \mathbf{Z} by z_{ij} throughout this supplement, and the corresponding partial derivatives of the above mappings by $\frac{\partial}{\partial z_{ij}}$.

Lemma E.1. *For multi-task elastic-net (i.e., $\tau > 0$ in (6)), the mapping $\mathbf{Z} \mapsto D^{-1}\mathbf{F}/\sqrt{n}$ is $n^{-\frac{1}{2}}L$ -Lipschitz with $L = 8 \max(1, (2\tau')^{-1})$, where $\tau' = \tau/\|\mathbf{\Sigma}\|_{\text{op}}$.*

Lemma E.2. *For multi-task group Lasso (i.e., $\tau = 0$ in (6)), we have*

- (1) *In $U_1 \cap U_2$, the map $\mathbf{Z} \mapsto D^{-1}\mathbf{F}/\sqrt{n}$ is $n^{-\frac{1}{2}}L$ -Lipschitz with $L = 8 \max(1, (2\eta)^{-1})$.*
- (2) *In $U_1 \cap U_2 \cap U_3$, the map $\mathbf{Z} \mapsto D^{-1}\mathbf{Z}^\top \mathbf{F}/n$ is $n^{-1/2}(1 + (2 + \sqrt{p/n})L)$ -Lipschitz, where $L = 8 \max(1, (2\eta)^{-1})$ as in (1).*

Corollary E.3. *Suppose that Assumption 4 holds, then*

- (1) *Under Assumption 4(i) that $\tau > 0$ and $\tau' = \tau/\|\mathbf{\Sigma}\|_{\text{op}}$, we have*

$$\sum_{ij} \left(\frac{\partial D}{\partial z_{ij}} \right)^2 \leq n^{-1} D^2 [4 \max(1, (2\tau')^{-1})]^2.$$

This implies that $nD^{-2} \sum_{ij} (\frac{\partial D}{\partial z_{ij}})^2 \leq C(\tau')$.

- (2) *Under Assumption 4(ii) that $\tau = 0$ and $\mathbb{P}(U_1) \rightarrow 1$, in the event $U_1 \cap U_2$, we have*

$$\sum_{ij} \left(\frac{\partial D}{\partial z_{ij}} \right)^2 \leq n^{-1} D^2 [4 \max(1, (2\eta)^{-1})]^2.$$

This implies that $nD^{-2} \sum_{ij} (\frac{\partial D}{\partial z_{ij}})^2 \leq C(\eta) = C(\gamma, c)$ since η is a constant that only depends on γ, c .

Note that with a fixed noise \mathbf{E} , Lemmas E.1 and E.2 guarantee that the map $\mathbf{Z} \mapsto \mathbf{F}$ is Lipschitz, hence the derivative exists almost everywhere by Rademacher's theorem. We present the formula for derivative of this map in Lemma E.4.

Lemma E.4. Recall $\mathbf{F} = \mathbf{Y} - \mathbf{X}\widehat{\mathbf{B}}$ with $\widehat{\mathbf{B}}$ defined in (6). Under Assumption 4(i) $\tau > 0$, or under Assumption 4(ii) $\tau = 0$ and in the event $U_1 \cap U_2$, for each $i, l \in [n], j \in [p], t \in [T]$, the following derivative exists almost everywhere and has the expression

$$\frac{\partial F_{lt}}{\partial z_{ij}} = D_{ij}^{lt} + \Delta_{ij}^{lt},$$

where $D_{ij}^{lt} = -(\mathbf{e}_j^\top \mathbf{H} \otimes \mathbf{e}_i^\top)(\mathbf{I}_{nT} - \mathbf{N})(\mathbf{e}_t \otimes \mathbf{e}_l)$, and $\Delta_{ij}^{lt} = -(\mathbf{e}_t^\top \otimes \mathbf{e}_l^\top)(\mathbf{I}_T \otimes \mathbf{X})\mathbf{M}^\dagger(\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}})(\mathbf{F}^\top \otimes \mathbf{I}_p)(\mathbf{e}_i \otimes \mathbf{e}_j)$. Furthermore, a straightforward calculation yields

$$\sum_{i=1}^n D_{ij}^{it} = -\mathbf{e}_j^\top \mathbf{H}(n\mathbf{I}_T - \widehat{\mathbf{A}})\mathbf{e}_t.$$

Lemma E.5. Suppose that Assumption 4 holds.

(1) Under Assumption 4(i) that $\tau > 0$ and $\tau' = \tau / \|\Sigma\|_{\text{op}}$, we have

$$\frac{1}{n} \sum_{ij} \left\| \frac{\partial(\mathbf{F}/D)}{\partial z_{ij}} \right\|_{\text{F}}^2 \leq \underbrace{4 \max(1, (\tau')^{-1}(T \wedge \frac{p}{n})) + 2n^{-1} [4 \max(1, (2\tau')^{-1})]^2}_{f(\tau', T, n, p)}.$$

(2) Under Assumption 4(ii) that $\tau = 0$ and $\mathbb{P}(U_1) \rightarrow 1$, in the event $U_1 \cap U_2$, we have

$$\frac{1}{n} \sum_{ij} \left\| \frac{\partial(\mathbf{F}/D)}{\partial z_{ij}} \right\|_{\text{F}}^2 \leq \underbrace{4 \max(1, (\eta)^{-1}(T \wedge \frac{p}{n})) + 2n^{-1} [4 \max(1, (2\eta)^{-1})]^2}_{f(\eta, T, n, p)}.$$

Furthermore, the right-hand side in (1) can be bounded from above by $C(\tau')(T \wedge \frac{p}{n})$, and the right-hand side in (2) can be bounded from above by $C(\gamma, c)$ in the regime $p/n \leq \gamma$.

E.2. Lipschitz and differential properties for a given, fixed design matrix

We also need to study Lipschitz and derivative properties of functions of the noise \mathbf{E} when the design \mathbf{X} is fixed. Formally, for a given and fixed design matrix \mathbf{X} , define the function $\mathbf{E} \mapsto \mathbf{F}(\mathbf{E})$ by the value $\mathbf{Y} - \mathbf{X}\widehat{\mathbf{B}}$ of the residual matrix when the observed data (\mathbf{X}, \mathbf{Y}) is $(\mathbf{X}, \mathbf{X}\mathbf{B}^* + \mathbf{E})$ and with $\widehat{\mathbf{B}}$ the estimator (6). Note that this map is 1-Lipschitz by (Bellec and Tsybakov, 2016, Proposition 3), Rademacher's theorem thus guarantees this map is differentiable almost everywhere. We denote its partial derivative by $\frac{\partial}{\partial E_{it}}$ for each entry $(E_{it})_{i \in [n], t \in [T]}$ of the noise matrix \mathbf{E} . We present its derivative formula in Lemma E.6 below.

Lemma E.6. *For each $i, l \in [n], t, t' \in [T]$, the following derivative exists almost everywhere and has the expression*

$$\frac{\partial F_{lt}}{\partial E_{it'}} = \mathbf{e}_l^\top \mathbf{e}_i \mathbf{e}_t^\top \mathbf{e}_{t'} - \mathbf{e}_l^\top (\mathbf{e}_t^\top \otimes \mathbf{X}) \mathbf{M}^\dagger (\mathbf{e}_{t'} \otimes \mathbf{X}^\top) \mathbf{e}_i.$$

As a consequence, we further have

$$\sum_{i=1}^n \frac{\partial F_{it}}{\partial E_{it'}} = \mathbf{e}_t^\top (n\mathbf{I}_T - \widehat{\mathbf{A}}) \mathbf{e}_{t'}, \quad \sum_{i=1}^n \frac{\partial \mathbf{e}_i^\top \mathbf{Z} \mathbf{H} \mathbf{e}_t}{\partial E_{it'}} = \mathbf{e}_t^\top \widehat{\mathbf{A}} \mathbf{e}_{t'}.$$

E.3. Probabilistic tools

We first list some useful variants of Stein's formulae and Gaussian-Poincaré inequalities. Let f' denote the derivative of a differentiable univariate function. For a differentiable vector-valued function $\mathbf{f}(\mathbf{z}) : \mathbb{R}^n \rightarrow \mathbb{R}^n$, denote its Jacobian (derivative) and divergence respectively by $\nabla \mathbf{f}(\mathbf{z})$ and $\text{div } \mathbf{f}(\mathbf{z})$, i.e., $[\nabla \mathbf{f}(\mathbf{z})]_{i,l} = \frac{\partial f_i(\mathbf{z})}{\partial z_l}$ for $i, l \in [n]$, and $\text{div } \mathbf{f}(\mathbf{z}) = \text{Tr}(\nabla \mathbf{f}(\mathbf{z}))$.

Lemma E.7 (Second-order Stein's formula Bellec and Zhang (2021)). *The following identities hold provided the involved derivatives exist a.e. and the expectations are finite.*

i) $z \sim \mathcal{N}(0, 1)$, $f : \mathbb{R} \rightarrow \mathbb{R}$, then

$$\mathbb{E}[(zf(z) - f'(z))^2] = \mathbb{E}[f(z)^2] + \mathbb{E}[(f'(z))^2].$$

ii) $\mathbf{z} \sim \mathcal{N}_n(\mathbf{0}, \mathbf{I}_n)$, $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$, then

$$\mathbb{E}[(\mathbf{z}^\top \mathbf{f}(\mathbf{z}) - \text{div } \mathbf{f}(\mathbf{z}))^2] = \mathbb{E}[\|\mathbf{f}(\mathbf{z})\|^2 + \text{Tr}[(\nabla \mathbf{f}(\mathbf{z}))^2]] \leq \mathbb{E}[\|\mathbf{f}(\mathbf{z})\|^2 + \|\nabla \mathbf{f}(\mathbf{z})\|_{\mathbb{F}}^2],$$

where the inequality uses Cauchy-Schwarz inequality.

iii) More generally, for $\mathbf{z} \sim \mathcal{N}_n(\mathbf{0}, \Sigma)$, $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$, then

$$\begin{aligned} \mathbb{E}[(\mathbf{z}^\top \mathbf{f}(\mathbf{z}) - \text{Tr}(\Sigma \nabla \mathbf{f}(\mathbf{z})))^2] &= \mathbb{E}[\|\Sigma^{\frac{1}{2}} \mathbf{f}(\mathbf{z})\|^2 + \text{Tr}[(\Sigma \nabla \mathbf{f}(\mathbf{z}))^2]] \\ &\leq \mathbb{E}[\|\Sigma^{\frac{1}{2}} \mathbf{f}(\mathbf{z})\|^2 + \|\Sigma \nabla \mathbf{f}(\mathbf{z})\|_{\mathbb{F}}^2], \end{aligned}$$

where the inequality uses Cauchy-Schwarz inequality.

Lemma E.8 (Gaussian-Poincaré inequality Boucheron, Lugosi and Massart (2013)). *The following inequalities hold provided the right-hand side derivatives exist a.e. and the expectations are finite.*

i) $z \sim \mathcal{N}(0, 1)$, $f : \mathbb{R} \rightarrow \mathbb{R}$, then

$$\text{Var}[f(z)] \leq \mathbb{E}[(f'(z))^2].$$

ii) $\mathbf{z} \sim \mathcal{N}_n(\mathbf{0}, \mathbf{I}_n)$, $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}$, then

$$\text{Var}[f(\mathbf{z})] \leq \mathbb{E}[\|\nabla f(\mathbf{z})\|^2].$$

iii) $\mathbf{z} \sim \mathcal{N}_n(\mathbf{0}, \mathbf{I}_n)$, $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, then

$$\mathbb{E}[\|\mathbf{f}(\mathbf{z}) - \mathbb{E}[\mathbf{f}(\mathbf{z})]\|^2] \leq \mathbb{E}[\|\nabla \mathbf{f}(\mathbf{z})\|_{\mathbb{F}}^2].$$

iv) More generally, for $\mathbf{z} \sim \mathcal{N}_n(\mathbf{0}, \Sigma)$, $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, then

$$\mathbb{E}[\|\mathbf{f}(\mathbf{z}) - \mathbb{E}[\mathbf{f}(\mathbf{z})]\|^2] \leq \mathbb{E}[\|\Sigma^{\frac{1}{2}} \nabla \mathbf{f}(\mathbf{z})\|_{\text{F}}^2].$$

Now we present a few important lemmas, whose proofs are based on Lemma E.7 and Lemma E.8.

Lemma E.9. Assume that Assumption 1 holds. For fixed \mathbf{X} , we have

$$\mathbb{E}\left[\|\mathbf{E}^\top \mathbf{F} / \tilde{D} - \mathbf{S}(n\mathbf{I}_T - \widehat{\mathbf{A}}) / \tilde{D}\|_{\text{F}}^2\right] \leq 4 \text{Tr}(\mathbf{S}),$$

where $\tilde{D} = (\|\mathbf{F}\|_{\text{F}}^2 + n \text{Tr}(\mathbf{S}))^{\frac{1}{2}}$.

Lemma E.10. Let $\mathbf{U}, \mathbf{V} : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^{n \times T}$ be two locally Lipschitz functions of \mathbf{Z} with i.i.d. $\mathcal{N}(0, 1)$ entries, then

$$\begin{aligned} & \mathbb{E}\left[\left\|\mathbf{U}^\top \mathbf{Z} \mathbf{V} - \sum_{j=1}^p \sum_{i=1}^n \frac{\partial}{\partial z_{ij}} (\mathbf{U}^\top \mathbf{e}_i \mathbf{e}_j^\top \mathbf{V})\right\|_{\text{F}}^2\right] \\ & \leq \mathbb{E}\|\mathbf{U}\|_{\text{F}}^2 \|\mathbf{V}\|_{\text{F}}^2 + \mathbb{E} \sum_{ij} \left[2\|\mathbf{V}\|_{\text{F}}^2 \left\|\frac{\partial \mathbf{U}}{\partial z_{ij}}\right\|_{\text{F}}^2 + 2\|\mathbf{U}\|_{\text{F}}^2 \left\|\frac{\partial \mathbf{V}}{\partial z_{ij}}\right\|_{\text{F}}^2\right]. \end{aligned}$$

Corollary E.11. Assume the same setting as Lemma E.10. If on some open set $\Omega \subset \mathbb{R}^{n \times p}$ with $\mathbb{P}(\Omega^c) \leq C/T$ for some constant C , we have (i) \mathbf{U} is $n^{-1/2}L_1$ -Lipschitz and $\|\mathbf{U}\|_{\text{F}} \leq 1$, (ii) \mathbf{V} is $n^{-1/2}L_2$ -Lipschitz and $\|\mathbf{V}\|_{\text{F}} \leq K$. Then

$$\begin{aligned} & \mathbb{E}\left[I(\Omega) \left\|\mathbf{U}^\top \mathbf{Z} \mathbf{V} - \sum_{j=1}^p \sum_{i=1}^n \frac{\partial}{\partial z_{ij}} (\mathbf{U}^\top \mathbf{e}_i \mathbf{e}_j^\top \mathbf{V})\right\|_{\text{F}}^2\right] \\ & \leq K^2 + 2C(K^2 L_1^2 + L_2^2) + 2\mathbb{E}\left[I(\Omega) \sum_{ij} \left(K^2 \left\|\frac{\partial \mathbf{U}}{\partial z_{ij}}\right\|_{\text{F}}^2 + \left\|\frac{\partial \mathbf{V}}{\partial z_{ij}}\right\|_{\text{F}}^2\right)\right]. \end{aligned}$$

Lemma E.12. Let $\mathbf{U}, \mathbf{V} : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^{n \times T}$ be two locally Lipschitz functions of \mathbf{Z} with i.i.d. $\mathcal{N}(0, 1)$ entries. Assume also that $\|\mathbf{U}\|_{\text{F}} \vee \|\mathbf{V}\|_{\text{F}} \leq 1$ almost surely. Then

$$\begin{aligned} & \mathbb{E}\left[\left\|p\mathbf{U}^\top \mathbf{V} - \sum_{j=1}^p \left(\sum_{i=1}^n \partial_{ij} \mathbf{U}^\top \mathbf{e}_i - \mathbf{U}^\top \mathbf{Z} \mathbf{e}_j\right) \left(\sum_{i=1}^n \partial_{ij} \mathbf{e}_i^\top \mathbf{V} - \mathbf{e}_j^\top \mathbf{Z}^\top \mathbf{V}\right)\right\|_{\text{F}}\right] \\ & \leq 2\|\mathbf{U}\|_{\partial} \|\mathbf{V}\|_{\partial} + \sqrt{p}(\sqrt{2} + (3 + \sqrt{2})(\|\mathbf{U}\|_{\partial} + \|\mathbf{V}\|_{\partial})), \end{aligned}$$

where $\partial_{ij} \mathbf{U} \stackrel{\text{def}}{=} \partial \mathbf{U} / \partial z_{ij}$, and $\|\mathbf{U}\|_{\partial} \stackrel{\text{def}}{=} \mathbb{E}[\sum_{i=1}^n \sum_{j=1}^p \|\partial_{ij} \mathbf{U}\|_{\text{F}}^2]^{\frac{1}{2}}$.

Lemma E.13. Let $\mathbf{N} = (\mathbf{I}_T \otimes \mathbf{X}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \mathbf{X}^\top)$, we have $\|\mathbf{N}\|_{\text{op}} \leq 1$.

Lemma E.14. *We have*

$$\begin{aligned} \sum_{ij} \mathbf{F}^\top \mathbf{Z} \mathbf{e}_j \mathbf{e}_i^\top \frac{\partial \mathbf{F}}{\partial z_{ij}} &= \mathbf{J}_1 - \mathbf{F}^\top \mathbf{Z} \mathbf{H} (n\mathbf{I}_T - \widehat{\mathbf{A}}), \\ \sum_{ij} \left(\frac{\partial \mathbf{F}}{\partial z_{ij}} \right)^\top \mathbf{Z} \mathbf{e}_j \mathbf{e}_i^\top \mathbf{F} &= \mathbf{J}_2 - \widehat{\mathbf{A}} \mathbf{F}^\top \mathbf{F}. \end{aligned}$$

where $\mathbf{J}_1 = \sum_{ij} \mathbf{F}^\top \mathbf{Z} \mathbf{e}_j \sum_t \Delta_{ij}^{it} \mathbf{e}_t^\top$ with $\|\mathbf{J}_1\|_F \leq n^{\frac{1}{2}} \|\mathbf{F}\|_F^2$, and $\mathbf{J}_2 = \sum_{ijlt} \mathbf{e}_t D_{ij}^{lt} \mathbf{e}_l^\top \mathbf{Z} \mathbf{e}_j \mathbf{e}_i^\top \mathbf{F}$ with $\|\mathbf{J}_2\|_F \leq n^{\frac{1}{2}} \|\mathbf{Z}\|_{\text{op}} \|\mathbf{H}\|_F \|\mathbf{F}\|_F$.

F. Proofs of Propositions B.3 to B.5

F.1. Proof of Proposition B.3

Proof of Proposition B.3. Note that $\mathbf{Q}_1 = \frac{\mathbf{E}^\top \mathbf{F} / \tilde{D} - \mathbf{S} (n\mathbf{I}_T - \widehat{\mathbf{A}}) / \tilde{D}}{\|\mathbf{S}\|_F^{\frac{1}{2}}}$, where $\tilde{D} = (\|\mathbf{F}\|_F^2 + n \text{Tr}(\mathbf{S}))^{\frac{1}{2}}$ is defined in Lemma E.9. Now, apply Lemma E.9, we obtain

$$\mathbb{E} \|\mathbf{Q}_1\|_F^2 = \mathbb{E} \left[\|\mathbf{E}^\top \mathbf{F} / \tilde{D} - \mathbf{S} (n\mathbf{I}_T - \widehat{\mathbf{A}}) / \tilde{D}\|_F^2 \right] \frac{1}{\text{Tr}(\mathbf{S})} \leq 4.$$

■

F.2. Proof of Proposition B.4

Proof of Proposition B.4. We first apply Lemma E.10. To be more specific, let $\mathbf{U} = n^{-\frac{1}{2}} \mathbf{F} / D$ and $\mathbf{V} = n^{-\frac{1}{2}} \mathbf{Z}^\top \mathbf{U}$ with $D = (\|\mathbf{F}\|_F^2 / n + \|\mathbf{H}\|_F^2)^{\frac{1}{2}}$, then $\|\mathbf{U}\|_F \leq 1$, $\|\mathbf{V}\|_F \leq n^{-\frac{1}{2}} \|\mathbf{Z}\|_{\text{op}}$. Lemma E.10 yields

$$\mathbb{E} \left(\left\| \mathbf{U}^\top \mathbf{Z} \mathbf{V} - \sum_{j=1}^p \sum_{i=1}^n \frac{\partial}{\partial z_{ij}} \left(\mathbf{U}^\top \mathbf{e}_i \mathbf{e}_j^\top \mathbf{V} \right) \right\|_F^2 \right) \quad (41)$$

$$\begin{aligned} &\leq \mathbb{E} [\|\mathbf{U}\|_F^2 \|\mathbf{V}\|_F^2] + \mathbb{E} \sum_{ij} \left[2 \|\mathbf{V}\|_F^2 \left\| \frac{\partial \mathbf{U}}{\partial z_{ij}} \right\|_F^2 + 2 \|\mathbf{U}\|_F^2 \left\| \frac{\partial \mathbf{V}}{\partial z_{ij}} \right\|_F^2 \right] \\ &\leq \mathbb{E} (n^{-1} \|\mathbf{Z}\|_{\text{op}}^2) + 2 \mathbb{E} \left(n^{-1} \|\mathbf{Z}\|_{\text{op}}^2 \sum_{ij} \left\| \frac{\partial \mathbf{U}}{\partial z_{ij}} \right\|_F^2 \right) + 2 \mathbb{E} \left(\sum_{ij} \left\| \frac{\partial \mathbf{V}}{\partial z_{ij}} \right\|_F^2 \right) \quad (42) \\ &\leq \frac{4p}{n} + \mathbb{E} \left(\frac{1}{n} \|\mathbf{Z}\|_{\text{op}}^2 \right) + 6 \mathbb{E} \left(n^{-1} \|\mathbf{Z}\|_{\text{op}}^2 \sum_{ij} \left\| \frac{\partial \mathbf{U}}{\partial z_{ij}} \right\|_F^2 \right), \end{aligned}$$

where the last inequality uses the following bound derived using $\mathbf{V} = n^{-\frac{1}{2}} \mathbf{Z}^\top \mathbf{U}$, and $\|\mathbf{U}\|_F \leq 1$,

$$\sum_{ij} \left\| \frac{\partial \mathbf{V}}{\partial z_{ij}} \right\|_F^2 = \sum_{ij} n^{-1} \left\| \frac{\partial \mathbf{Z}^\top}{\partial z_{ij}} \mathbf{U} + \mathbf{Z}^\top \frac{\partial \mathbf{U}}{\partial z_{ij}} \right\|_F^2$$

$$\begin{aligned}
&\leq 2n^{-1} \left(p \|\mathbf{U}\|_F^2 + \|\mathbf{Z}\|_{\text{op}}^2 \sum_{ij} \left\| \frac{\partial \mathbf{U}}{\partial z_{ij}} \right\|_F^2 \right) \\
&\leq \frac{2p}{n} + 2n^{-1} \|\mathbf{Z}\|_{\text{op}}^2 \sum_{ij} \left\| \frac{\partial \mathbf{U}}{\partial z_{ij}} \right\|_F^2.
\end{aligned} \tag{43}$$

Note that $\mathbb{E}[(n^{-\frac{1}{2}} \|\mathbf{Z}\|_{\text{op}})^2] \leq (1 + \sqrt{p/n})^2 + 1/n$ by (36). Now we establish the connection between \mathbf{Q}_2 and the term inside Frobenius norm in (41). By definitions of \mathbf{U} and \mathbf{V} ,

$$\mathbf{U}^\top \mathbf{Z} \mathbf{V} = n^{-\frac{3}{2}} D^{-2} \mathbf{F}^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{F}. \tag{44}$$

Next, by product rule,

$$\begin{aligned}
&\sum_{j=1}^p \sum_{i=1}^n \frac{\partial}{\partial z_{ij}} (\mathbf{U}^\top \mathbf{e}_i \mathbf{e}_j^\top \mathbf{V}) \\
&= n^{-\frac{3}{2}} \sum_{j=1}^p \sum_{i=1}^n \frac{\partial}{\partial z_{ij}} (\mathbf{F}^\top \mathbf{e}_i \mathbf{e}_j^\top \mathbf{Z}^\top \mathbf{F} D^{-2}) \\
&= n^{-\frac{3}{2}} \sum_{j=1}^p \sum_{i=1}^n \left(\underbrace{\frac{\partial \mathbf{F}^\top}{\partial z_{ij}} \mathbf{e}_i \mathbf{e}_j^\top \mathbf{Z}^\top \mathbf{F} D^{-2}}_{(i)} + \underbrace{\mathbf{F}^\top \mathbf{e}_i \mathbf{e}_j^\top \frac{\partial \mathbf{Z}^\top \mathbf{F}}{\partial z_{ij}} D^{-2}}_{(ii)} + \underbrace{\mathbf{F}^\top \mathbf{e}_i \mathbf{e}_j^\top \mathbf{Z}^\top \mathbf{F} \frac{\partial D^{-2}}{\partial z_{ij}}}_{(iii)} \right).
\end{aligned}$$

We now rewrite the above three terms (i), (ii) and (iii).

(i) For term (i), by Lemma E.14,

$$\begin{aligned}
&n^{-\frac{3}{2}} D^{-2} \sum_{j=1}^p \sum_{i=1}^n \frac{\partial \mathbf{F}^\top}{\partial z_{ij}} \mathbf{e}_i \mathbf{e}_j^\top \mathbf{Z}^\top \mathbf{F} \\
&= n^{-\frac{3}{2}} D^{-2} [\mathbf{J}_1 - \mathbf{F}^\top \mathbf{Z} \mathbf{H} (n\mathbf{I}_T - \widehat{\mathbf{A}})]^\top \\
&= n^{-\frac{3}{2}} D^{-2} [\mathbf{J}_1^\top - (n\mathbf{I}_T - \widehat{\mathbf{A}}) \mathbf{H}^\top \mathbf{Z}^\top \mathbf{F}].
\end{aligned}$$

(ii) For term (ii), by product rule and Lemma E.14,

$$\begin{aligned}
&n^{-\frac{3}{2}} D^{-2} \sum_{j=1}^p \sum_{i=1}^n \mathbf{F}^\top \mathbf{e}_i \mathbf{e}_j^\top \frac{\partial \mathbf{Z}^\top \mathbf{F}}{\partial z_{ij}} \\
&= n^{-\frac{3}{2}} D^{-2} \left(p \mathbf{F}^\top \mathbf{F} + \sum_{j=1}^p \sum_{i=1}^n \mathbf{F}^\top \mathbf{e}_i \mathbf{e}_j^\top \mathbf{Z}^\top \frac{\partial \mathbf{F}}{\partial z_{ij}} \right) \\
&= n^{-\frac{3}{2}} D^{-2} \left(p \mathbf{F}^\top \mathbf{F} + (\mathbf{J}_2 - \widehat{\mathbf{A}} \mathbf{F}^\top \mathbf{F})^\top \right) \\
&= n^{-\frac{3}{2}} D^{-2} (p \mathbf{F}^\top \mathbf{F} - \mathbf{F}^\top \mathbf{F} \widehat{\mathbf{A}} + \mathbf{J}_2^\top).
\end{aligned}$$

(iii) For term (iii), by chain rule,

$$\begin{aligned}
& n^{-\frac{3}{2}} \sum_{j=1}^p \sum_{i=1}^n \mathbf{F}^\top \mathbf{e}_i \mathbf{e}_j^\top \mathbf{Z}^\top \mathbf{F} \frac{\partial D^{-2}}{\partial z_{ij}} \\
&= -2n^{-\frac{3}{2}} D^{-3} \sum_{j=1}^p \sum_{i=1}^n \mathbf{F}^\top \mathbf{e}_i \mathbf{e}_j^\top \mathbf{Z}^\top \mathbf{F} \frac{\partial D}{\partial z_{ij}} \\
&= -n^{-\frac{3}{2}} D^{-2} \underbrace{\left(2D^{-1} \sum_{j=1}^p \sum_{i=1}^n \mathbf{F}^\top \mathbf{e}_i \mathbf{e}_j^\top \mathbf{Z}^\top \mathbf{F} \frac{\partial D}{\partial z_{ij}} \right)}_{\mathbf{J}_3}
\end{aligned}$$

Combining (44) and the above three expressions for (i)-(ii)-(iii),

$$\begin{aligned}
& \mathbf{U}^\top \mathbf{Z} \mathbf{V} - \sum_{j=1}^p \sum_{i=1}^n \frac{\partial}{\partial z_{ij}} \left(\mathbf{U}^\top \mathbf{e}_i \mathbf{e}_j^\top \mathbf{V} \right) \\
&= n^{-\frac{3}{2}} D^{-2} \left[\mathbf{F}^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{F} + (n\mathbf{I}_T - \widehat{\mathbf{A}}) \mathbf{H}^\top \mathbf{Z}^\top \mathbf{F} - p \mathbf{F}^\top \mathbf{F} + \mathbf{F}^\top \mathbf{F} \widehat{\mathbf{A}} \right. \\
&\quad \left. - \mathbf{J}_1^\top - \mathbf{J}_2^\top + \mathbf{J}_3 \right] \\
&= \mathbf{Q}_2 - n^{-\frac{3}{2}} D^{-2} (\mathbf{J}_1^\top + \mathbf{J}_2^\top - \mathbf{J}_3).
\end{aligned}$$

That is,

$$\mathbf{Q}_2 = \mathbf{U}^\top \mathbf{Z} \mathbf{V} - \sum_{j=1}^p \sum_{i=1}^n \frac{\partial}{\partial z_{ij}} \left(\mathbf{U}^\top \mathbf{e}_i \mathbf{e}_j^\top \mathbf{V} \right) + n^{-\frac{3}{2}} D^{-2} (\mathbf{J}_1^\top + \mathbf{J}_2^\top - \mathbf{J}_3). \quad (45)$$

Note that Lemma E.14 implies that

$$n^{-\frac{3}{2}} D^{-2} \|\mathbf{J}_1\|_{\mathbb{F}} \leq \frac{\|\mathbf{F}\|_{\mathbb{F}}^2/n}{\|\mathbf{F}\|_{\mathbb{F}}^2/n + \|\mathbf{H}\|_{\mathbb{F}}^2} \leq 1, \quad (46)$$

and

$$n^{-\frac{3}{2}} D^{-2} \|\mathbf{J}_2\|_{\mathbb{F}} \leq \left[n^{-\frac{1}{2}} \|\mathbf{Z}\|_{\text{op}} \frac{\|\mathbf{F}\|_{\mathbb{F}} n^{-\frac{1}{2}} \|\mathbf{H}\|_{\mathbb{F}}}{\|\mathbf{F}\|_{\mathbb{F}}^2/n + \|\mathbf{H}\|_{\mathbb{F}}^2} \right] \leq \frac{1}{2} (n^{-\frac{1}{2}} \|\mathbf{Z}\|_{\text{op}}). \quad (47)$$

Since $\mathbf{J}_3 = 2D^{-1} \sum_{j=1}^p \sum_{i=1}^n \mathbf{F}^\top \mathbf{e}_i \mathbf{e}_j^\top \mathbf{Z}^\top \mathbf{F} \frac{\partial D}{\partial z_{ij}}$, by Cauchy-Schwarz inequality

$$\begin{aligned}
n^{-3} D^{-4} \|\mathbf{J}_3\|_{\mathbb{F}}^2 &= \sum_{t,t'} ([\mathbf{J}_3]_{t,t'})^2 \\
&= 4n^{-3} D^{-6} \sum_{t,t'} \left[\sum_{i,j} \mathbf{e}_t^\top \mathbf{F}^\top \mathbf{e}_i \mathbf{e}_j^\top \mathbf{Z}^\top \mathbf{F} \mathbf{e}_{t'} \frac{\partial D}{\partial z_{ij}} \right]^2 \\
&\leq 4n^{-3} D^{-6} \sum_{t,t'} \left[\sum_{i,j} [\mathbf{e}_t^\top \mathbf{F}^\top \mathbf{e}_i \mathbf{e}_j^\top \mathbf{Z}^\top \mathbf{F} \mathbf{e}_{t'}]^2 \sum_{i,j} \left(\frac{\partial D}{\partial z_{ij}} \right)^2 \right]
\end{aligned}$$

$$\begin{aligned}
&= 4n^{-3}D^{-6}\|\mathbf{F}\|_{\mathbb{F}}^2\|\mathbf{Z}^\top\mathbf{F}\|_{\mathbb{F}}^2\sum_{i,j}\left(\frac{\partial D}{\partial z_{ij}}\right)^2 \\
&= 4n^{-3}D^{-6}\|\mathbf{F}\|_{\mathbb{F}}^4\|\mathbf{Z}\|_{\text{op}}^2\sum_{i,j}\left(\frac{\partial D}{\partial z_{ij}}\right)^2.
\end{aligned}$$

(1) Under Assumption 4(i) that $\tau > 0$. By Corollary E.3, we have $\sum_{i,j}\left(\frac{\partial D}{\partial z_{ij}}\right)^2 \leq C(\tau')n^{-1}D^2$. Then,

$$n^{-3}D^{-4}\|\mathbf{J}_3\|_{\mathbb{F}}^2 \leq C(\tau')n^{-4}D^2\|\mathbf{F}\|_{\mathbb{F}}^4\|\mathbf{Z}\|_{\text{op}}^2 \leq C(\tau')n^{-2}\|\mathbf{Z}\|_{\text{op}}^2, \quad (48)$$

by $\|\mathbf{F}\|_{\mathbb{F}}^2/(nD^2) \leq 1$.

By (45) and triangular inequality, we have

$$\begin{aligned}
\mathbb{E}[\|\mathbf{Q}_2\|_{\mathbb{F}}^2] &\leq 2\mathbb{E}\left[\|\mathbf{U}^\top\mathbf{Z}\mathbf{V} - \sum_{j=1}^p\sum_{i=1}^n\frac{\partial}{\partial z_{ij}}\left(\mathbf{U}^\top\mathbf{e}_i\mathbf{e}_j^\top\mathbf{V}\right)\|_{\mathbb{F}}^2\right] \\
&\quad + 2\mathbb{E}[\|n^{-\frac{3}{2}}D^{-2}(\mathbf{J}_1^\top + \mathbf{J}_2^\top - \mathbf{J}_3)\|_{\mathbb{F}}^2].
\end{aligned} \quad (49)$$

By (46), (47), (48), the second term in (49) can be upper bounded by

$$6\mathbb{E}\left[\left(1 + \frac{1}{4}n^{-1}\|\mathbf{Z}\|_{\text{op}}^2\right) + C(\tau')n^{-2}\|\mathbf{Z}\|_{\text{op}}^2\right] \leq C(\tau')\left(1 + \frac{p}{n}\right), \quad (50)$$

where the last inequality uses (36).

For the first term in (49), since $\sum_{ij}\|\frac{\partial \mathbf{U}}{\partial z_{ij}}\|_{\mathbb{F}}^2 \leq C(\tau')(T \wedge \frac{p}{n})$ by Lemma E.5, we have

$$\begin{aligned}
&\mathbb{E}\left(\left\|\mathbf{U}^\top\mathbf{Z}\mathbf{V} - \sum_{j=1}^p\sum_{i=1}^n\frac{\partial}{\partial z_{ij}}\left(\mathbf{U}^\top\mathbf{e}_i\mathbf{e}_j^\top\mathbf{V}\right)\right\|_{\mathbb{F}}^2\right) \\
&\leq \frac{4p}{n} + \mathbb{E}(n^{-1}\|\mathbf{Z}\|_{\text{op}}^2) + 6\mathbb{E}\left(n^{-1}\|\mathbf{Z}\|_{\text{op}}^2\sum_{ij}\left\|\frac{\partial \mathbf{U}}{\partial z_{ij}}\right\|_{\mathbb{F}}^2\right) \\
&\leq \frac{4p}{n} + [1 + C(\tau')(T \wedge \frac{p}{n})]\mathbb{E}(n^{-1}\|\mathbf{Z}\|_{\text{op}}^2) \\
&\leq \frac{4p}{n} + [1 + C(\tau')(T \wedge \frac{p}{n})][(1 + \sqrt{p/n})^2 + 1/n] \\
&\leq C(\tau')(T \wedge (1 + \frac{p}{n}))(1 + \frac{p}{n}).
\end{aligned}$$

Therefore, under Assumption 4(i), we obtain

$$\mathbb{E}[\|\mathbf{Q}_2\|_{\mathbb{F}}^2] \leq C(\tau')(T \wedge (1 + \frac{p}{n}))(1 + \frac{p}{n}).$$

(2) Under Assumption 4(ii), let $\Omega = U_1 \cap U_2 \cap U_3$, then we have $\mathbb{P}(\Omega^c) \leq C(\gamma, c)\frac{1}{T}$ by (37). By Lemma E.2, on Ω , we have (i) The map $\mathbf{Z} \mapsto \mathbf{U}$ is $n^{-1/2}L_1$ -Lipschitz, where $L_1 = 8\max(1, (2\eta)^{-1})$, and $\|\mathbf{U}\|_{\mathbb{F}} \leq 1$, (ii) The map $\mathbf{Z} \mapsto \mathbf{V}$ is $n^{-1/2}L_2$ -Lipschitz, where $L_2 = (1 + (2 + \sqrt{p/n})L_1)$, and

$\|V\|_F \leq (2 + \sqrt{p/n})$. Applying Corollary E.11 with $K = 2 + \sqrt{p/n}$ yields

$$\begin{aligned} & \mathbb{E} \left[I(\Omega) \left\| U^\top ZV - \sum_{j=1}^p \sum_{i=1}^n \frac{\partial}{\partial z_{ij}} (U^\top e_i e_j^\top V) \right\|_F^2 \right] \\ & \leq K^2 + C(\gamma, c)(K^2 L_1^2 + L_2^2) + 2\mathbb{E} \left[I(\Omega) \sum_{ij} \left(K^2 \left\| \frac{\partial U}{\partial z_{ij}} \right\|_F^2 + \left\| \frac{\partial V}{\partial z_{ij}} \right\|_F^2 \right) \right] \\ & \leq C(\gamma, c), \end{aligned}$$

where the last inequality holds because $K \leq C(\gamma)$, $L_1 = C(\gamma, c)$, $L_2 = C(\gamma, c)$, and on Ω , $\sum_{ij} \left\| \frac{\partial U}{\partial z_{ij}} \right\|_F^2 \leq C(\gamma, c)$ from Lemma E.5, and $\mathbb{E}[\left\| \frac{\partial V}{\partial z_{ij}} \right\|_F^2] \leq C(\gamma, c)$ by product rule. Therefore, under Assumption 4(i), we obtain

$$\begin{aligned} \mathbb{E}[I(\Omega) \|Q_2\|_F^2] & \leq 2\mathbb{E} \left[\left\| U^\top ZV - \sum_{j=1}^p \sum_{i=1}^n \frac{\partial}{\partial z_{ij}} (U^\top e_i e_j^\top V) \right\|_F^2 \right] \\ & \quad + 2\mathbb{E}[\|n^{-\frac{3}{2}} D^{-2} (J_1^\top + J_2^\top - J_3)\|_F^2] \\ & \leq C(\gamma, c), \end{aligned}$$

where the last inequality used (46), (47), and that $n^{-3} \mathbb{E}[D^{-4} \|J_3\|_F^2] \leq C(\gamma, c)$ in analogy to (48).

■

F.3. Proof of Proposition B.5

Proof of Proposition B.5. We will apply Lemma E.12. Let $U = V = n^{-\frac{1}{2}} D^{-1} F$ with $D = (\|F\|_F^2/n + \|H\|_F^2)^{\frac{1}{2}}$, then $\|U\|_F = \|V\|_F \leq 1$. Let $W_0 = pU^\top U - \sum_{j=1}^p (\sum_{i=1}^n \frac{\partial U^\top e_i}{\partial z_{ij}} - U^\top Z e_j)(\sum_{i=1}^n \frac{\partial U^\top e_i}{\partial z_{ij}} - U^\top Z e_j)^\top$, then Lemma E.12 gives

$$\begin{aligned} \mathbb{E}[\|W_0\|_F] & \leq 2\|U\|_\partial \|V\|_\partial + \sqrt{p}(\sqrt{2} + (3 + \sqrt{2})(\|U\|_\partial + \|V\|_\partial)) \\ & = 2\|U\|_\partial^2 + \sqrt{p}(\sqrt{2} + 2(3 + \sqrt{2})\|U\|_\partial). \end{aligned}$$

We will prove under Assumption 4(i), and the proof under Assumption 4(ii) on set $\Omega = U_1 \cap U_2 \cap U_3$ follows from almost similar arguments with τ' replaced by η , which is a constant that depends only on γ, c .

By definition of $\|U\|_\partial$ and Lemma E.5, $\|U\|_\partial^2 = \sum_{ij} \mathbb{E} \left\| \frac{\partial U}{\partial z_{ij}} \right\|_F^2 \leq C(\gamma, \tau')$. Thus $\mathbb{E}[\|W_0\|_F] \leq C(\tau', \gamma) \sqrt{p}$.

Now we establish the connection between W_0 and Q_3 . Since $U = n^{-\frac{1}{2}} D^{-1} F$, by product rule,

$$\begin{aligned} & \sum_{i=1}^n \frac{\partial U^\top e_i}{\partial z_{ij}} - U^\top Z e_j \\ & = n^{-\frac{1}{2}} \left(\sum_{i=1}^n \frac{\partial D^{-1} F^\top e_i}{\partial z_{ij}} - D^{-1} F^\top Z e_j \right) \end{aligned}$$

$$\begin{aligned}
&= n^{-\frac{1}{2}} \left(\sum_{i=1}^n D^{-1} \frac{\partial \mathbf{F}^\top \mathbf{e}_i}{\partial z_{ij}} - D^{-1} \mathbf{F}^\top \mathbf{Z} \mathbf{e}_j \right) + n^{-\frac{1}{2}} \sum_{i=1}^n \mathbf{F}^\top \mathbf{e}_i \frac{\partial D^{-1}}{\partial z_{ij}} \\
&= n^{-\frac{1}{2}} D^{-1} \left(\sum_{i=1}^n \frac{\partial \mathbf{F}^\top \mathbf{e}_i}{\partial z_{ij}} - \mathbf{F}^\top \mathbf{Z} \mathbf{e}_j \right) - n^{-\frac{1}{2}} D^{-2} \sum_{i=1}^n \mathbf{F}^\top \mathbf{e}_i \frac{\partial D}{\partial z_{ij}}.
\end{aligned}$$

For the first term in the last display, we have by Lemma E.4

$$\sum_{i=1}^n \frac{\partial \mathbf{e}_i^\top \mathbf{F}}{\partial z_{ij}} = \sum_{i=1}^n \sum_{t=1}^T \frac{\partial \mathbf{e}_i^\top \mathbf{F} \mathbf{e}_t}{\partial z_{ij}} \mathbf{e}_t^\top = \sum_{it} (D_{ij}^{it} + \Delta_{ij}^{it}) \mathbf{e}_t^\top = -\mathbf{e}_j^\top \mathbf{H} (n\mathbf{I}_T - \widehat{\mathbf{A}}) + \sum_{it} \Delta_{ij}^{it} \mathbf{e}_t^\top.$$

Hence,

$$\begin{aligned}
&\sum_{i=1}^n \frac{\partial \mathbf{U}^\top \mathbf{e}_i}{\partial z_{ij}} - \mathbf{U}^\top \mathbf{Z} \mathbf{e}_j \\
&= n^{-\frac{1}{2}} D^{-1} \left[-(n\mathbf{I}_T - \widehat{\mathbf{A}}) \mathbf{H}^\top \mathbf{e}_j + \sum_{it} \Delta_{ij}^{it} \mathbf{e}_t - \mathbf{F}^\top \mathbf{Z} \mathbf{e}_j \right] - n^{-\frac{1}{2}} D^{-2} \sum_{i=1}^n \mathbf{F}^\top \mathbf{e}_i \frac{\partial D}{\partial z_{ij}} \\
&= -n^{-\frac{1}{2}} D^{-1} \left[(n\mathbf{I}_T - \widehat{\mathbf{A}}) \mathbf{H}^\top + \mathbf{F}^\top \mathbf{Z} \right] \mathbf{e}_j + n^{-\frac{1}{2}} D^{-1} \sum_{it} \Delta_{ij}^{it} \mathbf{e}_t - n^{-\frac{1}{2}} D^{-2} \sum_{i=1}^n \mathbf{F}^\top \mathbf{e}_i \frac{\partial D}{\partial z_{ij}}. \quad (51)
\end{aligned}$$

Let $\mathbf{W}_1 = -n^{-\frac{1}{2}} D^{-1} \left[(n\mathbf{I}_T - \widehat{\mathbf{A}}) \mathbf{H}^\top + \mathbf{F}^\top \mathbf{Z} \right]$ be the first term in (51). For the second term in (51), recall $\Delta_{ij}^{it} = -(\mathbf{e}_t^\top \otimes \mathbf{e}_i^\top)(\mathbf{I}_T \otimes \mathbf{X}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}}) (\mathbf{F}^\top \otimes \mathbf{I}_p) (\mathbf{e}_i \otimes \mathbf{e}_j)$ in Lemma E.4,

$$\begin{aligned}
&n^{-\frac{1}{2}} D^{-1} \sum_{it} \Delta_{ij}^{it} \mathbf{e}_t \\
&= -n^{-\frac{1}{2}} D^{-1} \sum_{it} \mathbf{e}_t (\mathbf{e}_t^\top \otimes \mathbf{e}_i^\top) (\mathbf{I}_T \otimes \mathbf{X}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}}) (\mathbf{F}^\top \mathbf{e}_i \otimes \mathbf{e}_j) \\
&= -n^{-\frac{1}{2}} D^{-1} \sum_i (\mathbf{I}_T \otimes \mathbf{e}_i^\top \mathbf{X}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}}) (\mathbf{F}^\top \mathbf{e}_i \otimes \mathbf{I}_p) \mathbf{e}_j \\
&= -n^{-\frac{1}{2}} D^{-1} \sum_{i=1}^n (\mathbf{I}_T \otimes \mathbf{e}_i^\top \mathbf{X}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}}) (\mathbf{F}^\top \mathbf{e}_i \otimes \mathbf{I}_p) \mathbf{e}_j \\
&= \mathbf{W}_2 \mathbf{e}_j,
\end{aligned}$$

where $\mathbf{W}_2 = -n^{-\frac{1}{2}} D^{-1} \sum_{i=1}^n (\mathbf{I}_T \otimes \mathbf{e}_i^\top \mathbf{X}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}}) (\mathbf{F}^\top \mathbf{e}_i \otimes \mathbf{I}_p)$. For the third term in (51),

$$-n^{-\frac{1}{2}} D^{-2} \mathbf{F}^\top \sum_{i=1}^n \mathbf{e}_i \frac{\partial D}{\partial z_{ij}} = \mathbf{W}_3 \mathbf{e}_j,$$

where $\mathbf{W}_3 = -n^{-\frac{1}{2}} D^{-2} \mathbf{F}^\top \frac{\partial D}{\partial \mathbf{Z}}$, here we slightly abuse the notation and let $\frac{\partial D}{\partial \mathbf{Z}}$ denote the $n \times p$ matrix with (i, j) -th entry being $\frac{\partial D}{\partial z_{ij}}$. Therefore, (51) can be simplified as

$$\sum_{i=1}^n \frac{\partial \mathbf{U}^\top \mathbf{e}_i}{\partial z_{ij}} - \mathbf{U}^\top \mathbf{Z} \mathbf{e}_j = [\mathbf{W}_1 + \mathbf{W}_2 + \mathbf{W}_3] \mathbf{e}_j.$$

Furthermore,

$$\begin{aligned} \mathbf{W}_0 &= p \mathbf{U}^\top \mathbf{U} - \sum_{j=1}^p \left(\sum_{i=1}^n \frac{\partial \mathbf{U}^\top \mathbf{e}_i}{\partial z_{ij}} - \mathbf{U}^\top \mathbf{Z} \mathbf{e}_j \right) \left(\sum_{i=1}^n \frac{\partial \mathbf{U}^\top \mathbf{e}_i}{\partial z_{ij}} - \mathbf{U}^\top \mathbf{Z} \mathbf{e}_j \right)^\top \\ &= p \mathbf{U}^\top \mathbf{U} - [\mathbf{W}_1 + \mathbf{W}_2 + \mathbf{W}_3] [\mathbf{W}_1 + \mathbf{W}_2 + \mathbf{W}_3]^\top \\ &= n^{\frac{1}{2}} \mathbf{Q}_3 - \mathbf{W}_1 (\mathbf{W}_2 + \mathbf{W}_3)^\top - (\mathbf{W}_2 + \mathbf{W}_3) \mathbf{W}_1^\top - (\mathbf{W}_2 + \mathbf{W}_3) (\mathbf{W}_2 + \mathbf{W}_3)^\top, \end{aligned}$$

where the last equality is due to

$$\begin{aligned} & p \mathbf{U}^\top \mathbf{U} - \mathbf{W}_1 \mathbf{W}_1^\top \\ &= n^{-1} D^{-2} \left[p \mathbf{F}^\top \mathbf{F} - \mathbf{F}^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{F} - (n \mathbf{I}_T - \widehat{\mathbf{A}}) \mathbf{H}^\top \mathbf{H} (n \mathbf{I}_T - \widehat{\mathbf{A}}) \right. \\ &\quad \left. - (n \mathbf{I}_T - \widehat{\mathbf{A}}) \mathbf{H}^\top \mathbf{Z}^\top \mathbf{F} - \mathbf{F}^\top \mathbf{Z} \mathbf{H} (n \mathbf{I}_T - \widehat{\mathbf{A}}) \right] \\ &= n^{\frac{1}{2}} \mathbf{Q}_3. \end{aligned}$$

Therefore,

$$\mathbf{Q}_3 = n^{-\frac{1}{2}} \left[\mathbf{W}_0 + \mathbf{W}_1 (\mathbf{W}_2 + \mathbf{W}_3)^\top + (\mathbf{W}_2 + \mathbf{W}_3) \mathbf{W}_1^\top + (\mathbf{W}_2 + \mathbf{W}_3) (\mathbf{W}_2 + \mathbf{W}_3)^\top \right]. \quad (52)$$

We then bound the norms of $\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3$. For \mathbf{W}_1 ,

$$\begin{aligned} \|\mathbf{W}_1\|_F &= n^{-\frac{1}{2}} D^{-1} \|(n \mathbf{I}_T - \widehat{\mathbf{A}}) \mathbf{H}^\top + \mathbf{F}^\top \mathbf{Z}\|_F \\ &\leq n^{\frac{1}{2}} (D^{-1} \|\mathbf{H}\|_F + n^{-1} D^{-1} \|\mathbf{F}\|_F \|\mathbf{Z}\|_{\text{op}}) \\ &\leq n^{\frac{1}{2}} + n^{\frac{1}{2}} (D^{-1} \|\mathbf{F}\|_F / \sqrt{n} \|\mathbf{Z} / \sqrt{n}\|_{\text{op}}) \\ &\leq n^{\frac{1}{2}} (1 + \|\mathbf{Z} / \sqrt{n}\|_{\text{op}}), \end{aligned}$$

where we used $\|\mathbf{I}_T - \widehat{\mathbf{A}}/n\|_{\text{op}} \leq 1$ by Lemma B.1, $D^{-1} \|\mathbf{H}\|_F \leq 1$, and $D^{-1} \|\mathbf{F}\|_F / \sqrt{n} \leq 1$.

For \mathbf{W}_2 ,

$$\begin{aligned} \|\mathbf{W}_2\|_{\text{op}} &= n^{-\frac{1}{2}} D^{-1} \left\| \sum_{i=1}^n (\mathbf{I}_T \otimes \mathbf{e}_i^\top \mathbf{X}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}}) (\mathbf{F}^\top \mathbf{e}_i \otimes \mathbf{I}_p) \right\|_{\text{op}} \\ &\leq n^{-\frac{1}{2}} D^{-1} \sum_{i=1}^n \|(\mathbf{I}_T \otimes \mathbf{e}_i^\top \mathbf{X}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}}) (\mathbf{F}^\top \mathbf{e}_i \otimes \mathbf{I}_p)\|_{\text{op}} \end{aligned}$$

$$\begin{aligned}
&\leq n^{-\frac{1}{2}} D^{-1} \sum_{i=1}^n \|(I_T \otimes \mathbf{e}_i^\top X) \mathbf{M}^\dagger (I_T \otimes \Sigma^{\frac{1}{2}})\|_{\text{op}} \|(\mathbf{F}^\top \mathbf{e}_i \otimes I_p)\|_{\text{op}} \\
&\leq n^{-\frac{1}{2}} \|\Sigma\|_{\text{op}}^{\frac{1}{2}} D^{-1} \sum_{i=1}^n \|(I_T \otimes \mathbf{e}_i^\top X) (I_T \otimes (X_{\mathcal{J}}^\top X_{\mathcal{J}} + n\tau \mathbf{P}_{\mathcal{J}})^\dagger)\|_{\text{op}} \|(\mathbf{F}^\top \mathbf{e}_i \otimes I_p)\|_{\text{op}} \\
&= n^{-\frac{1}{2}} \|\Sigma\|_{\text{op}}^{\frac{1}{2}} D^{-1} \sum_{i=1}^n \|I_T \otimes [\mathbf{e}_i^\top X_{\mathcal{J}} (X_{\mathcal{J}}^\top X_{\mathcal{J}} + n\tau \mathbf{P}_{\mathcal{J}})^\dagger]\|_{\text{op}} \|\mathbf{F}^\top \mathbf{e}_i\| \\
&\leq n^{-\frac{1}{2}} \|\Sigma\|_{\text{op}}^{\frac{1}{2}} D^{-1} \|X_{\mathcal{J}} (X_{\mathcal{J}}^\top X_{\mathcal{J}} + n\tau \mathbf{P}_{\mathcal{J}})^\dagger\|_{\text{op}} \sum_{i=1}^n \|\mathbf{F}^\top \mathbf{e}_i\| \\
&\leq n^{-\frac{1}{2}} D^{-1} n^{-\frac{1}{2}} (\tau / \|\Sigma\|_{\text{op}})^{-\frac{1}{2}} n^{\frac{1}{2}} \|\mathbf{F}\|_{\text{F}} \\
&\leq (\tau')^{-\frac{1}{2}},
\end{aligned}$$

where the third inequality uses $\mathbf{M}^\dagger \preceq I_T \otimes (X_{\mathcal{J}}^\top X_{\mathcal{J}} + n\tau \mathbf{P}_{\mathcal{J}})^\dagger$, the fourth inequality uses the result that $\|\mathbf{e}_i^\top \mathbf{A}\|_{\text{op}} \leq \|\mathbf{A}\|_{\text{op}}$, the penultimate inequality uses $\|X_{\mathcal{J}} (X_{\mathcal{J}}^\top X_{\mathcal{J}} + n\tau \mathbf{P}_{\mathcal{J}})^\dagger\|_{\text{op}} \leq (n\tau)^{-1/2}$ and the Cauchy-Schwarz inequality, the last inequality follows from $n^{-1/2} D^{-1} \|\mathbf{F}\|_{\text{F}} \leq 1$. It immediately follows that $\|\mathbf{W}_2\|_{\text{F}} \leq \sqrt{T} C(\tau')$ since the rank of \mathbf{W}_2 is at most T .

For \mathbf{W}_3 , using $n^{-\frac{1}{2}} D^{-1} \|\mathbf{F}\|_{\text{F}} \leq 1$, and $\left\| \frac{\partial D}{\partial \mathbf{Z}} \right\|_{\text{F}} \leq n^{-\frac{1}{2}} D C(\tau')$ from Corollary E.3, we obtain

$$\|\mathbf{W}_3\|_{\text{F}} = n^{-\frac{1}{2}} D^{-2} \|\mathbf{F}^\top \frac{\partial D}{\partial \mathbf{Z}}\|_{\text{F}} \leq n^{-\frac{1}{2}} D^{-2} \|\mathbf{F}\|_{\text{F}} \left\| \frac{\partial D}{\partial \mathbf{Z}} \right\|_{\text{F}} \leq D^{-1} \left\| \frac{\partial D}{\partial \mathbf{Z}} \right\|_{\text{F}} \leq n^{-1/2} C(\tau').$$

The desired inequality follows by combining (52) and the bounds for $\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3$. ■

G. Proof of technical results in Section E

G.1. Proofs of results in Section E.1

Proof of Lemma E.1. Fixing E , if X, \bar{X} are two design matrices, and $\widehat{\mathbf{B}}, \bar{\mathbf{B}}$ are the two corresponding multi-task elastic net estimates. Let $\mathbf{Z} = X \Sigma^{-\frac{1}{2}}$, $\bar{\mathbf{Z}} = \bar{X} \Sigma^{-\frac{1}{2}}$, $\bar{\mathbf{H}} = \Sigma^{\frac{1}{2}} (\bar{\mathbf{B}} - \mathbf{B}^*)$, $\bar{\mathbf{F}} = \mathbf{Y} - \bar{X} \bar{\mathbf{B}}$, and $\bar{D} = [\|\bar{\mathbf{H}}\|_{\text{F}}^2 + \|\bar{\mathbf{F}}\|_{\text{F}}^2/n]^{\frac{1}{2}}$. Without loss of generality, we assume $\bar{D} \leq D$. Recall the multi-task elastic net estimate $\widehat{\mathbf{B}} = \arg \min_{\mathbf{B} \in \mathbb{R}^{p \times T}} (\frac{1}{2n} \|\mathbf{Y} - X \mathbf{B}\|_{\text{F}}^2 + g(\mathbf{B}))$, where $g(\mathbf{B}) = \lambda \|\mathbf{B}\|_{2,1} + \frac{\tau}{2} \|\mathbf{B}\|_{\text{F}}^2$. Define $\varphi : \mathbf{B} \mapsto \frac{1}{2n} \|\mathbf{E} + X(\mathbf{B}^* - \mathbf{B})\|_{\text{F}}^2 + g(\mathbf{B})$, $\psi : \mathbf{B} \mapsto \frac{1}{2n} \|X(\widehat{\mathbf{B}} - \mathbf{B})\|_{\text{F}}^2$ and $\zeta : \mathbf{B} \mapsto \varphi(\mathbf{B}) - \psi(\mathbf{B})$. When expanding the squares, it is clear that ζ is the sum of a linear function and a τ -strong convex penalty, thus ζ is τ -strongly convex of \mathbf{B} . Additivity of subdifferentials yields $\partial \varphi(\widehat{\mathbf{B}}) = \partial \zeta(\widehat{\mathbf{B}}) + \partial \psi(\widehat{\mathbf{B}}) = \partial \zeta(\widehat{\mathbf{B}})$. By optimality of $\widehat{\mathbf{B}}$ we have $\mathbf{0}_{p \times T} \in \partial \varphi(\widehat{\mathbf{B}})$, thus $\mathbf{0}_{p \times T} \in \partial \zeta(\widehat{\mathbf{B}})$. By strong convexity of ζ , $\zeta(\bar{\mathbf{B}}) - \zeta(\widehat{\mathbf{B}}) \geq \langle \partial \zeta(\widehat{\mathbf{B}}), \bar{\mathbf{B}} - \widehat{\mathbf{B}} \rangle + \frac{\tau}{2} \|\bar{\mathbf{B}} - \widehat{\mathbf{B}}\|_{\text{F}}^2 = \frac{\tau}{2} \|\bar{\mathbf{B}} - \widehat{\mathbf{B}}\|_{\text{F}}^2$, which can further be rewritten as

$$\|X(\widehat{\mathbf{B}} - \bar{\mathbf{B}})\|_{\text{F}}^2 + n\tau \|\bar{\mathbf{B}} - \widehat{\mathbf{B}}\|_{\text{F}}^2 \leq \|\mathbf{E} - X(\bar{\mathbf{B}} - \mathbf{B}^*)\|_{\text{F}}^2 - \|\mathbf{E} - X(\widehat{\mathbf{B}} - \mathbf{B}^*)\|_{\text{F}}^2 + 2n(g(\bar{\mathbf{B}}) - g(\widehat{\mathbf{B}})),$$

i.e.,

$$\|\mathbf{Z}(\mathbf{H} - \bar{\mathbf{H}})\|_{\text{F}}^2 + n\tau \|\Sigma^{-\frac{1}{2}}(\mathbf{H} - \bar{\mathbf{H}})\|_{\text{F}}^2 \leq \|\mathbf{E} - \mathbf{Z}\bar{\mathbf{H}}\|_{\text{F}}^2 - \|\mathbf{E} - \mathbf{Z}\mathbf{H}\|_{\text{F}}^2 + 2n(g(\bar{\mathbf{B}}) - g(\widehat{\mathbf{B}})).$$

Summing the above inequality with its counterpart obtained by replacing $(X, \hat{\mathbf{B}}, \mathbf{H})$ with $(\bar{X}, \bar{\mathbf{B}}, \bar{\mathbf{H}})$, we have

$$\begin{aligned}
& (LHS) \\
& \stackrel{\text{def}}{=} \|\mathbf{Z}(\mathbf{H} - \bar{\mathbf{H}})\|_{\text{F}}^2 + \|\bar{\mathbf{Z}}(\mathbf{H} - \bar{\mathbf{H}})\|_{\text{F}}^2 + 2n\tau' \|\mathbf{H} - \bar{\mathbf{H}}\|_{\text{F}}^2 \\
& \leq \|\mathbf{Z}(\mathbf{H} - \bar{\mathbf{H}})\|_{\text{F}}^2 + \|\bar{\mathbf{Z}}(\mathbf{H} - \bar{\mathbf{H}})\|_{\text{F}}^2 + 2n\tau \|\boldsymbol{\Sigma}^{-\frac{1}{2}}(\mathbf{H} - \bar{\mathbf{H}})\|_{\text{F}}^2 \\
& \leq \|\mathbf{E} - \mathbf{Z}\bar{\mathbf{H}}\|_{\text{F}}^2 - \|\mathbf{E} - \mathbf{Z}\mathbf{H}\|_{\text{F}}^2 + \|\mathbf{E} - \bar{\mathbf{Z}}\mathbf{H}\|_{\text{F}}^2 - \|\mathbf{E} - \bar{\mathbf{Z}}\bar{\mathbf{H}}\|_{\text{F}}^2 \\
& = \langle \mathbf{Z}(\mathbf{H} - \bar{\mathbf{H}}), \mathbf{F} + \bar{\mathbf{F}} + (\bar{\mathbf{Z}} - \mathbf{Z})\bar{\mathbf{H}} \rangle + \langle -\bar{\mathbf{Z}}(\mathbf{H} - \bar{\mathbf{H}}), \mathbf{F} + \bar{\mathbf{F}} + (\mathbf{Z} - \bar{\mathbf{Z}})\mathbf{H} \rangle \\
& = \langle (\mathbf{Z} - \bar{\mathbf{Z}})(\mathbf{H} - \bar{\mathbf{H}}), \mathbf{F} + \bar{\mathbf{F}} \rangle + \langle \mathbf{Z}(\mathbf{H} - \bar{\mathbf{H}}), (\bar{\mathbf{Z}} - \mathbf{Z})\bar{\mathbf{H}} \rangle + \langle \bar{\mathbf{Z}}(\bar{\mathbf{H}} - \mathbf{H}), (\mathbf{Z} - \bar{\mathbf{Z}})\mathbf{H} \rangle \\
& \leq \|\mathbf{Z} - \bar{\mathbf{Z}}\|_{\text{op}} \|\mathbf{H} - \bar{\mathbf{H}}\|_{\text{F}} (\|\mathbf{F}\|_{\text{F}} + \|\bar{\mathbf{F}}\|_{\text{F}}) + \|\mathbf{Z} - \bar{\mathbf{Z}}\|_{\text{op}} \|\mathbf{Z}(\mathbf{H} - \bar{\mathbf{H}})\|_{\text{F}} \|\bar{\mathbf{H}}\|_{\text{F}} \\
& \quad + \|\mathbf{Z} - \bar{\mathbf{Z}}\|_{\text{op}} \|\bar{\mathbf{Z}}(\mathbf{H} - \bar{\mathbf{H}})\|_{\text{F}} \|\mathbf{H}\|_{\text{F}} \\
& \leq \|\mathbf{Z} - \bar{\mathbf{Z}}\|_{\text{op}} \left[\sqrt{\frac{(LHS)}{2n\tau'}} (\|\mathbf{F}\|_{\text{F}} + \|\bar{\mathbf{F}}\|_{\text{F}}) + \sqrt{(LHS)} (\|\bar{\mathbf{H}}\|_{\text{F}} + \|\mathbf{H}\|_{\text{F}}) \right] \\
& \leq \|\mathbf{Z} - \bar{\mathbf{Z}}\|_{\text{op}} \sqrt{(LHS)} (D + \bar{D}) \max(1, (2\tau')^{-\frac{1}{2}})
\end{aligned}$$

where $\tau' = \tau \phi_{\min}(\boldsymbol{\Sigma}^{-1}) = \tau / \|\boldsymbol{\Sigma}\|_{\text{op}}$. That is,

$$\sqrt{(LHS)} \leq \|\mathbf{Z} - \bar{\mathbf{Z}}\|_{\text{op}} 2D \max(1, (2\tau')^{-\frac{1}{2}}).$$

Therefore,

$$\begin{aligned}
& n^{-\frac{1}{2}} \|\mathbf{F} - \bar{\mathbf{F}}\|_{\text{F}} = n^{-\frac{1}{2}} \|\mathbf{Z}\mathbf{H} - \bar{\mathbf{Z}}\bar{\mathbf{H}}\|_{\text{F}} \\
& \leq n^{-\frac{1}{2}} [\|\mathbf{Z}(\mathbf{H} - \bar{\mathbf{H}})\|_{\text{F}} + \|(\mathbf{Z} - \bar{\mathbf{Z}})\bar{\mathbf{H}}\|_{\text{F}}] \\
& \leq n^{-\frac{1}{2}} [\|\mathbf{Z}(\mathbf{H} - \bar{\mathbf{H}})\|_{\text{F}} + \|\mathbf{Z} - \bar{\mathbf{Z}}\|_{\text{op}} \|\bar{\mathbf{H}}\|_{\text{F}}] \\
& \leq n^{-\frac{1}{2}} [\sqrt{(LHS)} + \|\mathbf{Z} - \bar{\mathbf{Z}}\|_{\text{op}} D] \\
& \leq n^{-\frac{1}{2}} \|\mathbf{Z} - \bar{\mathbf{Z}}\|_{\text{op}} D [2 \max(1, (2\tau')^{-\frac{1}{2}}) + 1].
\end{aligned}$$

So far we obtained

$$\begin{aligned}
& \|\mathbf{H} - \bar{\mathbf{H}}\|_{\text{F}} \leq \sqrt{\frac{(LHS)}{2n\tau'}} \leq n^{-\frac{1}{2}} \|\mathbf{Z} - \bar{\mathbf{Z}}\|_{\text{op}} D (2\tau')^{-\frac{1}{2}} 2 \max(1, (2\tau')^{-\frac{1}{2}}), \\
& n^{-\frac{1}{2}} \|\mathbf{F} - \bar{\mathbf{F}}\|_{\text{F}} \leq n^{-\frac{1}{2}} \|\mathbf{Z} - \bar{\mathbf{Z}}\|_{\text{op}} D [2 \max(1, (2\tau')^{-\frac{1}{2}}) + 1].
\end{aligned}$$

Let $\mathbf{Q} = [\mathbf{H}^{\top}, \mathbf{F}^{\top} / \sqrt{n}]^{\top}$ and $\bar{\mathbf{Q}} = [\bar{\mathbf{H}}^{\top}, \bar{\mathbf{F}}^{\top} / \sqrt{n}]^{\top}$, then $D = \|\mathbf{Q}\|_{\text{F}}$, $\bar{D} = \|\bar{\mathbf{Q}}\|_{\text{F}}$. By triangular inequality,

$$\begin{aligned}
|D - \bar{D}| & \leq \|\mathbf{Q} - \bar{\mathbf{Q}}\|_{\text{F}} \leq \|\mathbf{H} - \bar{\mathbf{H}}\|_{\text{F}} + \|\mathbf{F} - \bar{\mathbf{F}}\|_{\text{F}} / \sqrt{n} \\
& \leq n^{-\frac{1}{2}} \|\mathbf{Z} - \bar{\mathbf{Z}}\|_{\text{op}} D [4 \max(1, (2\tau')^{-1})],
\end{aligned}$$

where the last inequality uses the elementary inequality $\max(a, b)(a + b) \leq 2[\max(a, b)]^2$ for $a, b > 0$ with $a = 1, b = (2\tau')^{-\frac{1}{2}}$. Let $\frac{\partial D}{\partial \mathbf{Z}} \stackrel{\text{def}}{=} \frac{\partial D}{\partial \text{vec}(\mathbf{Z})} \in \mathbb{R}^{1 \times np}$, then $\|\frac{\partial D}{\partial \mathbf{Z}}\| \leq n^{-\frac{1}{2}}DL_1$ with $L_1 = [4\max(1, (2\tau')^{-1})]$. Hence,

$$\sum_{ij} \left(\frac{\partial D}{\partial z_{ij}} \right)^2 = \left\| \frac{\partial D}{\partial \mathbf{Z}} \right\|^2 \leq n^{-1}D^2L_1^2.$$

Furthermore, by triangle inequality

$$\begin{aligned} \left\| \frac{\mathbf{Q}}{D} - \frac{\bar{\mathbf{Q}}}{\bar{D}} \right\|_F &\leq \frac{1}{D} \|\mathbf{Q} - \bar{\mathbf{Q}}\|_F + \left| \frac{1}{D} - \frac{1}{\bar{D}} \right| \|\bar{\mathbf{Q}}\|_F \\ &= \frac{1}{D} \|\mathbf{Q} - \bar{\mathbf{Q}}\|_F + \frac{|D - \bar{D}|}{D\bar{D}} \|\bar{\mathbf{Q}}\|_F \\ &\leq \frac{1}{D} \|\mathbf{Q} - \bar{\mathbf{Q}}\|_F + \frac{1}{D} \|\mathbf{Q} - \bar{\mathbf{Q}}\|_F \\ &\leq n^{-\frac{1}{2}} \|\mathbf{Z} - \bar{\mathbf{Z}}\|_{\text{op}} L, \end{aligned}$$

where $L = 8\max(1, (2\tau')^{-1})$. Therefore, when $\tau > 0$, we obtain the two mappings $\mathbf{Z} \mapsto D^{-1}\mathbf{F}/\sqrt{n}$, and $\mathbf{Z} \mapsto D^{-1}\mathbf{H}$ are both $n^{-\frac{1}{2}}L$ -lipschitz with $L = 8\max(1, (2\tau')^{-1})$, where $\tau' = \tau/\|\Sigma\|_{\text{op}}$. ■

The proof of Lemma E.2 uses a similar argument as proof of Lemma E.1, we present it here for completeness.

Proof of Lemma E.2. For multi-task group Lasso ($\tau = 0$), we restrict our analysis in the event $U_1 \cap U_2$, where $U_1 = \{\|\widehat{\mathbf{B}}\|_0 \leq n(1-c)/2\}$, $U_2 = \{\inf_{\mathbf{b} \in \mathbb{R}^p: \|\mathbf{b}\|_0 \leq (1-c)n} \|\mathbf{X}\mathbf{b}\|^2 / (n\|\Sigma^{\frac{1}{2}}\mathbf{b}\|^2) > \eta\}$.

Since the only randomness of the problem comes from \mathbf{X} and \mathbf{E} , there exists a measurable set \mathcal{U} such that $U_1 \cap U_2 = \{(\mathbf{X}, \mathbf{E}) \in \mathcal{U}\}$. For some noise matrix \mathbf{E} , consider $\mathbf{X}, \bar{\mathbf{X}}$ two design matrices such that $(\mathbf{X}, \mathbf{E}) \in \mathcal{U}$ and $(\bar{\mathbf{X}}, \mathbf{E}) \in \mathcal{U}$. We slightly abuse the notation and let $\bar{\mathbf{B}}, \bar{\mathbf{B}}$ denote the two corresponding multi-task group-Lasso estimates. Thus, the row sparsity of $\bar{\mathbf{B}} - \bar{\mathbf{B}}$ is at most $n(1-c)$. Let $\bar{\mathbf{H}} = \bar{\mathbf{B}} - \mathbf{B}^*$, $\bar{\mathbf{F}} = \mathbf{Y} - \bar{\mathbf{X}}\bar{\mathbf{B}}$, and $\bar{D} = [\|\bar{\mathbf{H}}\|_F^2 + \|\bar{\mathbf{F}}\|_F^2/n]^{\frac{1}{2}}$. Without loss of generality, we assume $\bar{D} \leq D$. Since when $\tau = 0$, the multi-task group Lasso estimate is $\widehat{\mathbf{B}} = \arg \min_{\mathbf{B} \in \mathbb{R}^{p \times T}} (\frac{1}{2n} \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_F^2 + g(\mathbf{B}))$, where $g(\mathbf{B}) = \lambda \|\mathbf{B}\|_{2,1}$. Define $\varphi: \mathbf{B} \mapsto \frac{1}{2n} \|\mathbf{E} + \mathbf{X}(\mathbf{B}^* - \mathbf{B})\|_F^2 + g(\mathbf{B})$, $\psi: \mathbf{B} \mapsto \frac{1}{2n} \|\mathbf{X}(\widehat{\mathbf{B}} - \mathbf{B})\|_F^2$ and $\zeta: \mathbf{B} \mapsto \varphi(\mathbf{B}) - \psi(\mathbf{B})$. Under $\tau = 0$, by the same arguments in proof of E.1 with the same functions $\varphi(\cdot), \psi(\cdot), \zeta(\cdot)$, we obtain

$$\|\mathbf{X}(\widehat{\mathbf{B}} - \bar{\mathbf{B}})\|_F^2 \leq \|\mathbf{E} - \mathbf{X}(\bar{\mathbf{B}} - \mathbf{B}^*)\|_F^2 - \|\mathbf{E} - \mathbf{X}(\widehat{\mathbf{B}} - \mathbf{B}^*)\|_F^2 + 2n(g(\bar{\mathbf{B}}) - g(\widehat{\mathbf{B}})).$$

Summing the above inequality with its counterpart obtained by replacing $(\mathbf{X}, \widehat{\mathbf{B}}, \mathbf{H})$ with $(\bar{\mathbf{X}}, \bar{\mathbf{B}}, \bar{\mathbf{H}})$, we have

$$\begin{aligned} &\|\mathbf{X}(\widehat{\mathbf{B}} - \bar{\mathbf{B}})\|_F^2 + \|\bar{\mathbf{X}}(\bar{\mathbf{B}} - \bar{\mathbf{B}})\|_F^2 \\ &\leq \|\mathbf{E} - \mathbf{Z}\bar{\mathbf{H}}\|_F^2 - \|\mathbf{E} - \mathbf{Z}\mathbf{H}\|_F^2 + \|\mathbf{E} - \bar{\mathbf{Z}}\mathbf{H}\|_F^2 - \|\mathbf{E} - \bar{\mathbf{Z}}\bar{\mathbf{H}}\|_F^2. \end{aligned}$$

Note that in event $U_1 \cap U_2$, we have

$$\eta n \|\Sigma^{\frac{1}{2}}(\widehat{\mathbf{B}} - \bar{\mathbf{B}})\|_F^2 \leq \|\mathbf{X}(\widehat{\mathbf{B}} - \bar{\mathbf{B}})\|_F^2, \quad \eta n \|\Sigma^{\frac{1}{2}}(\bar{\mathbf{B}} - \bar{\mathbf{B}})\|_F^2 \leq \|\bar{\mathbf{X}}(\bar{\mathbf{B}} - \bar{\mathbf{B}})\|_F^2.$$

Thus, $2\eta n \|(\hat{\mathbf{H}} - \bar{\mathbf{H}})\|_{\mathbb{F}}^2 \leq \|\mathbf{Z}(\mathbf{H} - \bar{\mathbf{H}})\|_{\mathbb{F}}^2 + \|\bar{\mathbf{Z}}(\mathbf{H} - \bar{\mathbf{H}})\|_{\mathbb{F}}^2$, and

$$\begin{aligned} (LHS) &\stackrel{\text{def}}{=} \max(2\eta n \|\mathbf{H} - \bar{\mathbf{H}}\|_{\mathbb{F}}^2, \|\mathbf{Z}(\mathbf{H} - \bar{\mathbf{H}})\|_{\mathbb{F}}^2 + \|\bar{\mathbf{Z}}(\mathbf{H} - \bar{\mathbf{H}})\|_{\mathbb{F}}^2) \\ &= \|\mathbf{Z}(\mathbf{H} - \bar{\mathbf{H}})\|_{\mathbb{F}}^2 + \|\bar{\mathbf{Z}}(\mathbf{H} - \bar{\mathbf{H}})\|_{\mathbb{F}}^2 \\ &\leq \|\mathbf{E} - \mathbf{Z}\bar{\mathbf{H}}\|_{\mathbb{F}}^2 - \|\mathbf{E} - \mathbf{Z}\mathbf{H}\|_{\mathbb{F}}^2 + \|\mathbf{E} - \bar{\mathbf{Z}}\mathbf{H}\|_{\mathbb{F}}^2 - \|\mathbf{E} - \bar{\mathbf{Z}}\bar{\mathbf{H}}\|_{\mathbb{F}}^2. \end{aligned}$$

Now, in $U_1 \cap U_2$, the Lipschitz property of the map $\mathbf{Z} \mapsto D^{-1}\mathbf{F}/\sqrt{n}$ follows from the same arguments in proof of Lemma E.1, with τ' in E.1 replaced by η in this proof.

Furthermore, in the event $U_1 \cap U_2 \cap U_3$, the Lipschitz property of $\mathbf{Z} \mapsto D^{-1}\mathbf{Z}^\top \mathbf{F}/n$ follows by triangle inequality. To see this, let $\mathbf{U} = D^{-1}\mathbf{F}/\sqrt{n}$, and $\mathbf{V} = D^{-1}\mathbf{Z}^\top \mathbf{F}/n = n^{-1/2}\mathbf{Z}^\top \mathbf{U}$, thus by triangle inequality

$$\begin{aligned} \|\mathbf{V} - \bar{\mathbf{V}}\|_{\text{op}} &= n^{-1/2} \|\mathbf{Z}^\top \mathbf{U} - \bar{\mathbf{Z}}^\top \bar{\mathbf{U}}\|_{\text{op}} \\ &= n^{-1/2} [\|(\mathbf{Z} - \bar{\mathbf{Z}})^\top \mathbf{U}\|_{\text{op}} + \|\bar{\mathbf{Z}}^\top (\mathbf{U} - \bar{\mathbf{U}})\|_{\text{op}}] \\ &\leq n^{-1/2} [\|\mathbf{Z} - \bar{\mathbf{Z}}\|_{\text{op}} + \|\bar{\mathbf{Z}}\|_{\text{op}} \|\mathbf{U} - \bar{\mathbf{U}}\|_{\text{op}}] \\ &\leq n^{-1/2} (1 + n^{-1/2} \|\bar{\mathbf{Z}}\|_{\text{op}} L) \|\mathbf{Z} - \bar{\mathbf{Z}}\|_{\text{op}} \\ &\leq n^{-1/2} (1 + (2 + \sqrt{p/n})L). \end{aligned}$$

where the last line uses $\|\bar{\mathbf{Z}}\|_{\text{op}} \leq 2\sqrt{n} + \sqrt{p}$ in the event U_3 . ■

Proof of Corollary E.3. Corollary E.3 (1) is a direct consequence of the intermediate result $|D - \bar{D}| \leq n^{-\frac{1}{2}} \|\mathbf{Z} - \bar{\mathbf{Z}}\|_{\text{op}} D [4 \max(1, (2\tau')^{-1})]$ in proof of Lemma E.1, while Corollary E.3 (2) is a direct consequence of the intermediate result $|D - \bar{D}| \leq n^{-\frac{1}{2}} \|\mathbf{Z} - \bar{\mathbf{Z}}\|_{\text{op}} D [4 \max(1, (2\eta)^{-1})]$ in proof of Lemma E.2. ■

Before proving the derivative formula, we restate $\hat{\mathbf{B}}$ (defined in (6) of the full paper) below,

$$\hat{\mathbf{B}} = \arg \min_{\mathbf{B} \in \mathbb{R}^{p \times T}} \left(\frac{1}{2n} \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_{\mathbb{F}}^2 + \lambda \|\mathbf{B}\|_{2,1} + \frac{\tau}{2} \|\mathbf{B}\|_{\mathbb{F}}^2 \right), \quad (53)$$

where $\|\mathbf{B}\|_{2,1} = \sum_{j=1}^P \|\mathbf{B}^\top \mathbf{e}_j\|_2$.

For the reader's convenience, we recall some useful notations. $\mathbf{P}_{\hat{\mathcal{J}}} = \sum_{k \in \hat{\mathcal{J}}} \mathbf{e}_k \mathbf{e}_k^\top$. For each $k \in \hat{\mathcal{J}}$, $\mathbf{H}^{(k)} = \lambda \|\hat{\mathbf{B}}^\top \mathbf{e}_k\|_2^{-1} \left(\mathbf{I}_T - \hat{\mathbf{B}}^\top \mathbf{e}_k \mathbf{e}_k^\top \hat{\mathbf{B}} \right) \|\hat{\mathbf{B}}^\top \mathbf{e}_k\|_2^{-2}$. $\tilde{\mathbf{H}} = \sum_{k \in \hat{\mathcal{J}}} (\mathbf{H}^{(k)} \otimes \mathbf{e}_k \mathbf{e}_k^\top)$. $\mathbf{M}_1 = \mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + \tau n \mathbf{P}_{\hat{\mathcal{J}}})$, $\mathbf{M} = \mathbf{M}_1 + n \tilde{\mathbf{H}} \in \mathbb{R}^{pT \times pT}$, and let $\mathbf{N} = (\mathbf{I}_T \otimes \mathbf{X}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \mathbf{X}^\top)$.

Proof of Lemma E.4. We first derive $\frac{\partial F_{lt}}{\partial x_{ij}}$. Since $\mathbf{F} = \mathbf{Y} - \mathbf{X}\hat{\mathbf{B}} = \mathbf{E} - \mathbf{X}(\hat{\mathbf{B}} - \mathbf{B}^*)$, by product rule,

$$\frac{\partial F_{lt}}{\partial x_{ij}} = \mathbf{e}_l^\top \frac{\partial \mathbf{E} - \mathbf{X}(\hat{\mathbf{B}} - \mathbf{B}^*)}{\partial x_{ij}} \mathbf{e}_t = -\mathbf{e}_l^\top (\dot{\mathbf{X}}(\hat{\mathbf{B}} - \mathbf{B}^*) + \mathbf{X}\dot{\mathbf{B}}) \mathbf{e}_t,$$

where $\dot{\mathbf{X}} \stackrel{\text{def}}{=} \frac{\partial \mathbf{X}}{\partial x_{ij}} = \mathbf{e}_i \mathbf{e}_j^\top$, and $\dot{\mathbf{B}} \stackrel{\text{def}}{=} \frac{\partial \hat{\mathbf{B}}}{\partial x_{ij}}$.

Now we derive $\text{vec}(\dot{\mathbf{B}})$ from KKT conditions for $\hat{\mathbf{B}}$ defined in (53):

1) For $k \in \hat{\mathcal{J}}$, i.e., $\widehat{\mathbf{B}}^\top \mathbf{e}_k \neq \mathbf{0}$,

$$\mathbf{e}_k^\top \mathbf{X}^\top [\mathbf{E} - \mathbf{X}(\widehat{\mathbf{B}} - \mathbf{B}^*)] - n\tau \mathbf{e}_k^\top \widehat{\mathbf{B}} = \frac{n\lambda}{\|\widehat{\mathbf{B}}^\top \mathbf{e}_k\|_2} \mathbf{e}_k^\top \widehat{\mathbf{B}} \in \mathbb{R}^{1 \times T}.$$

2) For $k \notin \hat{\mathcal{J}}$, i.e., $\widehat{\mathbf{B}}^\top \mathbf{e}_k = \mathbf{0}$,

$$\left\| \mathbf{e}_k^\top \mathbf{X}^\top [\mathbf{E} - \mathbf{X}(\widehat{\mathbf{B}} - \mathbf{B}^*)] - n\tau \mathbf{e}_k^\top \widehat{\mathbf{B}} \right\| < n\lambda.$$

Here the strict inequality is guaranteed by Proposition 2.3 of Bellec (2020).

Keeping \mathbf{E} fixed, differentiation of the above display for $k \in \hat{\mathcal{J}}$ w.r.t. x_{ij} yields

$$\mathbf{e}_k^\top \left[\dot{\mathbf{X}}^\top \mathbf{F} - \mathbf{X}^\top [\dot{\mathbf{X}}(\widehat{\mathbf{B}} - \mathbf{B}^*) + \mathbf{X}\dot{\mathbf{B}}] - n\tau \dot{\mathbf{B}} \right] = n\mathbf{e}_k^\top \dot{\mathbf{B}} \mathbf{H}^{(k)},$$

with $\mathbf{H}^{(k)} = \lambda \|\widehat{\mathbf{B}}^\top \mathbf{e}_k\|_2^{-1} \left(\mathbf{I}_T - \widehat{\mathbf{B}}^\top \mathbf{e}_k \mathbf{e}_k^\top \widehat{\mathbf{B}} \|\widehat{\mathbf{B}}^\top \mathbf{e}_k\|_2^{-2} \right) \in \mathbb{R}^{T \times T}$. Rearranging and using $\dot{\mathbf{X}} = \mathbf{e}_i \mathbf{e}_j^\top$,

$$\mathbf{e}_k^\top \left[\mathbf{e}_j \mathbf{e}_i^\top \mathbf{F} - \mathbf{X}^\top \mathbf{e}_i \mathbf{e}_j^\top (\widehat{\mathbf{B}} - \mathbf{B}^*) \right] = \mathbf{e}_k^\top [(\mathbf{X}^\top \mathbf{X} + n\tau \mathbf{I}_p) \dot{\mathbf{B}} + n\dot{\mathbf{B}} \mathbf{H}^{(k)}].$$

Recall $\mathbf{P}_{\hat{\mathcal{J}}} = \sum_{k \in \hat{\mathcal{J}}} \mathbf{e}_k \mathbf{e}_k^\top$. Multiplying by \mathbf{e}_k to the left and summing over $k \in \hat{\mathcal{J}}$, we obtain

$$\mathbf{P}_{\hat{\mathcal{J}}} \left[\mathbf{e}_j \mathbf{e}_i^\top \mathbf{F} - \mathbf{X}^\top \mathbf{e}_i \mathbf{e}_j^\top (\widehat{\mathbf{B}} - \mathbf{B}^*) \right] = \mathbf{P}_{\hat{\mathcal{J}}} (\mathbf{X}^\top \mathbf{X} + n\tau \mathbf{I}_p) \dot{\mathbf{B}} + n \sum_{k \in \hat{\mathcal{J}}} \mathbf{e}_k \mathbf{e}_k^\top \dot{\mathbf{B}} \mathbf{H}^{(k)}.$$

Since $\hat{\mathcal{J}}$ is locally constant in a small neighborhood of \mathbf{X} , $\widehat{\mathbf{B}}_{\hat{\mathcal{J}}^c} = \mathbf{0}$, $\text{supp}(\dot{\mathbf{B}}) \subseteq \hat{\mathcal{J}}$. Hence, $\mathbf{P}_{\hat{\mathcal{J}}} \dot{\mathbf{B}} = \dot{\mathbf{B}}$, and $\mathbf{X} \dot{\mathbf{B}} = \mathbf{X}_{\hat{\mathcal{J}}} \dot{\mathbf{B}}$. The above display can be rewritten as

$$\mathbf{P}_{\hat{\mathcal{J}}} \mathbf{e}_j \mathbf{e}_i^\top \mathbf{F} - \mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{e}_i \mathbf{e}_j^\top (\widehat{\mathbf{B}} - \mathbf{B}^*) = (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}}) \dot{\mathbf{B}} + n \sum_{k \in \hat{\mathcal{J}}} \mathbf{e}_k \mathbf{e}_k^\top \dot{\mathbf{B}} \mathbf{H}^{(k)}.$$

Vectorizing the above display using property $\text{vec}(\mathbf{ABC}) = (\mathbf{C}^\top \otimes \mathbf{A}) \text{vec}(\mathbf{A})$ yields

$$\begin{aligned} & (\mathbf{F}^\top \otimes \mathbf{P}_{\hat{\mathcal{J}}} \mathbf{e}_j) \text{vec}(\mathbf{e}_i^\top) - ((\widehat{\mathbf{B}} - \mathbf{B}^*)^\top \mathbf{e}_j \otimes \mathbf{X}_{\hat{\mathcal{J}}}^\top) \text{vec}(\mathbf{e}_i) \\ &= [\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}}) + n \sum_{k \in \hat{\mathcal{J}}} (\mathbf{H}^{(k)} \otimes \mathbf{e}_k \mathbf{e}_k^\top)] \text{vec}(\dot{\mathbf{B}}) \\ &= (\mathbf{M}_1 + n\tilde{\mathbf{H}}) \text{vec}(\dot{\mathbf{B}}) \\ &= \mathbf{M} \text{vec}(\dot{\mathbf{B}}), \end{aligned}$$

where $\mathbf{M}_1 = \mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}})$, and $\tilde{\mathbf{H}} = \sum_{k \in \hat{\mathcal{J}}} (\mathbf{H}^{(k)} \otimes \mathbf{e}_k \mathbf{e}_k^\top)$.

Under Assumption 4(i) that $\tau > 0$, it's obviously that $\text{rank}(\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}})) = T|\hat{\mathcal{J}}|$. Under Assumption 4(ii) that $\tau = 0$ with $\mathbb{P}(U_1) \rightarrow 1$. In the event $U_1 \cap U_2$, we know $\text{rank}(\mathbf{X}_{\hat{\mathcal{J}}}) = |\hat{\mathcal{J}}|$ from (Bellec and Romon, 2021, Lemma C.4), hence $\text{rank}(\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}})) = T|\hat{\mathcal{J}}|$. In either of the above two scenarios, we thus have $\dim(\ker(\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}}))) = T(p - |\hat{\mathcal{J}}|)$ by rank-nullity theorem. Since $[\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}})](\mathbf{e}_t \otimes \mathbf{e}_k) = \mathbf{0}$ for $t \in [T], k \in \hat{\mathcal{J}}^c$. Let $V = \{(\mathbf{e}_t \otimes \mathbf{e}_k) : t \in$

$[T], k \in \hat{\mathcal{J}}^c\}$ be a vector space, then the elements of V are linear independent, and $\dim(V) = T(p - |\hat{\mathcal{J}}|)$. Thus, V forms a basis for $\ker(\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}}))$. Since for any $(\mathbf{e}_t \otimes \mathbf{e}_k) \in V$, we also have $\tilde{\mathbf{H}}(\mathbf{e}_t \otimes \mathbf{e}_k) = \mathbf{0}$, $\ker(\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}})) \subseteq \ker(\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}} + n\tilde{\mathbf{H}}))$. On the other hand, if any vector \mathbf{v} s.t. $[\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}}) + n\tilde{\mathbf{H}}]\mathbf{v} = \mathbf{0}$, since these matrices are all positive semi-definite, we have $\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}})\mathbf{v} = \mathbf{0}$, which implies that $\ker(\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}}) + n\tilde{\mathbf{H}}) \subseteq \ker(\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}}))$. Therefore,

$$\begin{aligned} \ker(\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}}) + n\tilde{\mathbf{H}}) &= \ker(\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}})) \\ &= \text{span}\{(\mathbf{e}_t \otimes \mathbf{e}_k) : t \in [T], k \in \hat{\mathcal{J}}^c\}, \end{aligned}$$

and

$$\text{range}(\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}}) + n\tilde{\mathbf{H}}) = \text{span}\{(\mathbf{e}_t \otimes \mathbf{e}_k) : t \in [T], k \in \hat{\mathcal{J}}\}.$$

Since $\dot{\mathbf{B}} = \mathbf{P}_{\hat{\mathcal{J}}} \dot{\mathbf{B}}$, $\text{vec}(\dot{\mathbf{B}}) = (\mathbf{I}_T \otimes \mathbf{P}_{\hat{\mathcal{J}}}) \text{vec}(\dot{\mathbf{B}})$, then $\text{vec}(\dot{\mathbf{B}}) \in \text{col}(\mathbf{I}_T \otimes \mathbf{P}_{\hat{\mathcal{J}}}) = \text{range}(\mathbf{M})$. Since \mathbf{M} is symmetric, $\mathbf{M}^\dagger \mathbf{M}$ is the orthogonal projection on the range of \mathbf{M} . Therefore,

$$\text{vec}(\dot{\mathbf{B}}) = \mathbf{M}^\dagger \mathbf{M} \text{vec}(\dot{\mathbf{B}}) = \mathbf{M}^\dagger [(\mathbf{F}^\top \otimes \mathbf{e}_j) - ((\widehat{\mathbf{B}} - \mathbf{B}^*)^\top \mathbf{e}_j \otimes \mathbf{X}^\top)] \mathbf{e}_i. \quad (54)$$

Since $\text{supp}(\dot{\mathbf{B}}) \subseteq \hat{\mathcal{J}}$, $\mathbf{X} \dot{\mathbf{B}} = \mathbf{X}_{\hat{\mathcal{J}}} \dot{\mathbf{B}}$, we have

$$\begin{aligned} \frac{\partial F_{lt}}{\partial x_{ij}} &= -\mathbf{e}_l^\top (\dot{\mathbf{X}}(\widehat{\mathbf{B}} - \mathbf{B}^*) + \mathbf{X} \dot{\mathbf{B}}) \mathbf{e}_t \\ &= -(\mathbf{e}_l^\top \mathbf{e}_i \mathbf{e}_j^\top (\widehat{\mathbf{B}} - \mathbf{B}^*) \mathbf{e}_t + \mathbf{e}_l^\top \mathbf{X}_{\hat{\mathcal{J}}} \dot{\mathbf{B}} \mathbf{e}_t) \\ &= -(\mathbf{e}_l^\top \mathbf{e}_i \mathbf{e}_j^\top (\widehat{\mathbf{B}} - \mathbf{B}^*) \mathbf{e}_t + (\mathbf{e}_l^\top \otimes \mathbf{e}_l^\top \mathbf{X}_{\hat{\mathcal{J}}}) \text{vec}(\dot{\mathbf{B}})) \\ &= -\mathbf{e}_l^\top \mathbf{e}_i \mathbf{e}_j^\top (\widehat{\mathbf{B}} - \mathbf{B}^*) \mathbf{e}_t - (\mathbf{e}_l^\top \otimes \mathbf{e}_l^\top \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger [(\mathbf{F}^\top \otimes \mathbf{e}_j) - ((\widehat{\mathbf{B}} - \mathbf{B}^*)^\top \mathbf{e}_j \otimes \mathbf{X}^\top)] \mathbf{e}_i \\ &= -(\mathbf{e}_j^\top (\widehat{\mathbf{B}} - \mathbf{B}^*) \otimes \mathbf{e}_l^\top) (\mathbf{e}_t \otimes \mathbf{e}_l) + (\mathbf{e}_l^\top \otimes \mathbf{e}_l^\top \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger ((\widehat{\mathbf{B}} - \mathbf{B}^*)^\top \mathbf{e}_j \otimes \mathbf{X}^\top \mathbf{e}_i) \\ &\quad - (\mathbf{e}_l^\top \otimes \mathbf{e}_l^\top \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger (\mathbf{F}^\top \otimes \mathbf{e}_j) \mathbf{e}_i \\ &= -(\mathbf{e}_j^\top (\widehat{\mathbf{B}} - \mathbf{B}^*) \otimes \mathbf{e}_l^\top) (\mathbf{e}_t \otimes \mathbf{e}_l) + (\mathbf{e}_j^\top (\widehat{\mathbf{B}} - \mathbf{B}^*) \otimes \mathbf{e}_l^\top) \mathbf{N} (\mathbf{e}_t \otimes \mathbf{e}_l) \\ &\quad - (\mathbf{e}_l^\top \otimes \mathbf{e}_l^\top \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger (\mathbf{F}^\top \otimes \mathbf{I}_p) (\mathbf{e}_i \otimes \mathbf{e}_j) \\ &= -(\mathbf{e}_j^\top (\widehat{\mathbf{B}} - \mathbf{B}^*) \otimes \mathbf{e}_l^\top) (\mathbf{I}_{nT} - \mathbf{N}) (\mathbf{e}_t \otimes \mathbf{e}_l) - (\mathbf{e}_l^\top \otimes \mathbf{e}_l^\top \mathbf{X}) \mathbf{M}^\dagger (\mathbf{F}^\top \otimes \mathbf{I}_p) (\mathbf{e}_i \otimes \mathbf{e}_j) \end{aligned}$$

Now we calculate $\frac{\partial F_{lt}}{\partial z_{ij}}$. Since $\mathbf{X} = \mathbf{Z} \boldsymbol{\Sigma}^{\frac{1}{2}}$, $x_{ik} = \sum_{j=1}^p z_{ij} (\boldsymbol{\Sigma}^{\frac{1}{2}})_{jk}$, $\frac{\partial x_{ik}}{\partial z_{ij}} = (\boldsymbol{\Sigma}^{\frac{1}{2}})_{jk}$,

$$\frac{\partial F_{lt}}{\partial z_{ij}} = \sum_{k=1}^p \frac{\partial F_{lt}}{\partial x_{ik}} \frac{\partial x_{ik}}{\partial z_{ij}} = \sum_{k=1}^p \frac{\partial F_{lt}}{\partial x_{ik}} (\boldsymbol{\Sigma}^{\frac{1}{2}})_{jk} = D_{ij}^{lt} + \Delta_{ij}^{lt},$$

where

$$D_{ij}^{lt} = - \sum_{k=1}^p (\mathbf{e}_k^\top (\widehat{\mathbf{B}} - \mathbf{B}^*) \otimes \mathbf{e}_l^\top) (\mathbf{I}_{nT} - \mathbf{N}) (\mathbf{e}_t \otimes \mathbf{e}_l) (\boldsymbol{\Sigma}^{\frac{1}{2}})_{jk}$$

$$\begin{aligned}
&= -(\mathbf{e}_j^\top \Sigma^{\frac{1}{2}} (\widehat{\mathbf{B}} - \mathbf{B}^*) \otimes \mathbf{e}_i^\top) (\mathbf{I}_{nT} - \mathbf{N}) (\mathbf{e}_t \otimes \mathbf{e}_l) \\
&= -(\mathbf{e}_j^\top \mathbf{H} \otimes \mathbf{e}_i^\top) (\mathbf{I}_{nT} - \mathbf{N}) (\mathbf{e}_t \otimes \mathbf{e}_l),
\end{aligned}$$

and

$$\begin{aligned}
\Delta_{ij}^{lt} &= -\sum_{k=1}^p (\mathbf{e}_t^\top \otimes \mathbf{e}_l^\top) (\mathbf{I}_T \otimes \mathbf{X}) \mathbf{M}^\dagger (\mathbf{F}^\top \otimes \mathbf{I}_p) (\mathbf{e}_i \otimes \mathbf{e}_k) (\Sigma^{\frac{1}{2}})_{jk} \\
&= -(\mathbf{e}_t^\top \otimes \mathbf{e}_l^\top) (\mathbf{I}_T \otimes \mathbf{X}) \mathbf{M}^\dagger (\mathbf{F}^\top \otimes \mathbf{I}_p) (\mathbf{e}_i \otimes \Sigma^{\frac{1}{2}} \mathbf{e}_j) \\
&= -(\mathbf{e}_t^\top \otimes \mathbf{e}_l^\top) (\mathbf{I}_T \otimes \mathbf{X}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}}) (\mathbf{F}^\top \otimes \mathbf{I}_p) (\mathbf{e}_i \otimes \mathbf{e}_j)
\end{aligned}$$

It follows that

$$\begin{aligned}
\sum_{i=1}^n D_{ij}^{it} &= -\sum_{i=1}^n (\mathbf{e}_j^\top \mathbf{H} \otimes \mathbf{e}_i^\top) (\mathbf{I}_{nT} - \mathbf{N}) (\mathbf{e}_t \otimes \mathbf{e}_i) \\
&= -\mathbf{e}_j^\top \mathbf{H} \left[\sum_{i=1}^n (\mathbf{I}_T \otimes \mathbf{e}_i^\top) (\mathbf{I}_{nT} - \mathbf{N}) (\mathbf{I}_T \otimes \mathbf{e}_i) \right] \mathbf{e}_t \\
&= -\mathbf{e}_j^\top \mathbf{H} (n\mathbf{I}_T - \widehat{\mathbf{A}}) \mathbf{e}_t,
\end{aligned}$$

where the last line follows from definition of $\widehat{\mathbf{A}}$ in (10). ■

Proof of Lemma E.5. (1) For $\tau > 0$, by formula of $\frac{\partial F_{lt}}{\partial z_{ij}}$ in Lemma E.4, we have

$$\begin{aligned}
\sum_{ij} \left\| \frac{\partial \mathbf{F}}{\partial z_{ij}} \right\|_{\mathbb{F}}^2 &= \sum_{ij} \sum_{lt} \left(\frac{\partial F_{lt}}{\partial z_{ij}} \right)^2 = \sum_{ij} \sum_{lt} (D_{ij}^{lt} + \Delta_{ij}^{lt})^2 \\
&\leq 2 \sum_{ij, lt} (D_{ij}^{lt})^2 + 2 \sum_{ij, lt} (\Delta_{ij}^{lt})^2 \\
&= 2 \|(\mathbf{H} \otimes \mathbf{I}_n) (\mathbf{I}_{nT} - \mathbf{N})\|_{\mathbb{F}}^2 + 2 \|(\mathbf{I}_T \otimes \mathbf{X}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}}) (\mathbf{F}^\top \otimes \mathbf{I}_p)\|_{\mathbb{F}}^2 \\
&\leq 2n \|\mathbf{H}\|_{\mathbb{F}}^2 + 2 \|(\mathbf{I}_T \otimes \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}}) (\mathbf{F}^\top \otimes \mathbf{I}_p)\|_{\mathbb{F}}^2.
\end{aligned}$$

Since $0 \preceq \mathbf{M}^\dagger \preceq \mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + \tau n \mathbf{P}_{\hat{\mathcal{J}}})^\dagger$,

$$\begin{aligned}
&\|(\mathbf{I}_T \otimes \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}}) (\mathbf{F}^\top \otimes \mathbf{I}_p)\|_{\mathbb{F}}^2 \\
&\leq \|(\mathbf{I}_T \otimes \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}})\|_{\text{op}}^2 \|(\mathbf{F}^\top \otimes \mathbf{I}_p)\|_{\mathbb{F}}^2 \\
&\leq p \|\Sigma\|_{\text{op}} \|\mathbf{F}\|_{\mathbb{F}}^2 \|(\mathbf{I}_T \otimes \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger\|_{\text{op}}^2 \\
&\leq p \|\Sigma\|_{\text{op}} \|\mathbf{F}\|_{\mathbb{F}}^2 \|(\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + \tau n \mathbf{P}_{\hat{\mathcal{J}}})^\dagger \mathbf{X}_{\hat{\mathcal{J}}}^\top\|_{\text{op}}^2 \\
&\leq \frac{p}{n\tau} \|\Sigma\|_{\text{op}} \|\mathbf{F}\|_{\mathbb{F}}^2
\end{aligned}$$

$$= \frac{p}{n\tau'} \|\mathbf{F}\|_{\mathbb{F}}^2,$$

where the last inequality uses $\|(\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + \tau n \mathbf{P}_{\hat{\mathcal{J}}})^\dagger \mathbf{X}_{\hat{\mathcal{J}}}^\top\|_{\text{op}} \leq (n\tau)^{-1}$.

On the other hand, we also have

$$\begin{aligned} & \|(\mathbf{I}_T \otimes \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \boldsymbol{\Sigma}^{\frac{1}{2}}) (\mathbf{F}^\top \otimes \mathbf{I}_p)\|_{\mathbb{F}}^2 \\ & \leq \|(\mathbf{I}_T \otimes \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger\|_{\mathbb{F}}^2 \|(\mathbf{I}_T \otimes \boldsymbol{\Sigma}^{\frac{1}{2}}) (\mathbf{F}^\top \otimes \mathbf{I}_p)\|_{\text{op}}^2 \\ & \leq \|(\mathbf{I}_T \otimes \mathbf{X}_{\hat{\mathcal{J}}}) (\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + \tau n \mathbf{P}_{\hat{\mathcal{J}}})^\dagger)\|_{\mathbb{F}}^2 \|\mathbf{F}\|_{\mathbb{F}}^2 \|\boldsymbol{\Sigma}\|_{\text{op}} \\ & \leq T \|\mathbf{X}_{\hat{\mathcal{J}}} (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + \tau n \mathbf{P}_{\hat{\mathcal{J}}})^\dagger\|_{\mathbb{F}}^2 \|\mathbf{F}\|_{\mathbb{F}}^2 \|\boldsymbol{\Sigma}\|_{\text{op}} \\ & \leq T \text{Tr} [(\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + \tau n \mathbf{P}_{\hat{\mathcal{J}}})^\dagger \mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + \tau n \mathbf{P}_{\hat{\mathcal{J}}})^\dagger] \|\mathbf{F}\|_{\mathbb{F}}^2 \|\boldsymbol{\Sigma}\|_{\text{op}} \\ & \leq T \text{Tr} [(\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + \tau n \mathbf{P}_{\hat{\mathcal{J}}})^\dagger] \|\mathbf{F}\|_{\mathbb{F}}^2 \|\boldsymbol{\Sigma}\|_{\text{op}} \\ & \leq T(\tau)^{-1} \|\mathbf{F}\|_{\mathbb{F}}^2 \|\boldsymbol{\Sigma}\|_{\text{op}} \\ & \leq T \text{Tr} [(\tau n \mathbf{P}_{\hat{\mathcal{J}}})^\dagger] \|\mathbf{F}\|_{\mathbb{F}}^2 \|\boldsymbol{\Sigma}\|_{\text{op}} \\ & \leq T(\tau)^{-1} \|\mathbf{F}\|_{\mathbb{F}}^2 \|\boldsymbol{\Sigma}\|_{\text{op}} \\ & = \frac{T}{\tau'} \|\mathbf{F}\|_{\mathbb{F}}^2, \end{aligned}$$

Therefore,

$$\begin{aligned} \frac{1}{n} \sum_{ij} \left\| \frac{\partial \mathbf{F}}{\partial z_{ij}} \right\|_{\mathbb{F}}^2 & \leq 2 \|\mathbf{H}\|_{\mathbb{F}}^2 + 2(\tau')^{-1} (T \wedge \frac{p}{n}) \|\mathbf{F}\|_{\mathbb{F}}^2 / n \\ & \leq 2 \max(1, (\tau')^{-1} (T \wedge \frac{p}{n})) (\|\mathbf{F}\|_{\mathbb{F}}^2 / n + \|\mathbf{H}\|_{\mathbb{F}}^2) \\ & = 2 \max(1, (\tau')^{-1} (T \wedge \frac{p}{n})) D^2. \end{aligned}$$

Now by product rule and triangle inequality

$$\begin{aligned} & \frac{1}{n} \sum_{ij} \left\| \frac{\partial \mathbf{F}/D}{\partial z_{ij}} \right\|_{\mathbb{F}}^2 \\ & \leq 2D^{-2} \frac{1}{n} \sum_{ij} \left\| \frac{\partial \mathbf{F}}{\partial z_{ij}} \right\|_{\mathbb{F}}^2 + 2 \frac{1}{n} \sum_{ij} \left\| \mathbf{F} \frac{\partial D^{-1}}{\partial z_{ij}} \right\|_{\mathbb{F}}^2 \\ & = 2D^{-2} \frac{1}{n} \sum_{ij} \left\| \frac{\partial \mathbf{F}}{\partial z_{ij}} \right\|_{\mathbb{F}}^2 + 2D^{-4} \frac{1}{n} \|\mathbf{F}\|_{\mathbb{F}}^2 \sum_{ij} \left(\frac{\partial D}{\partial z_{ij}} \right)^2 \\ & \leq 2D^{-2} \frac{1}{n} \sum_{ij} \left\| \frac{\partial \mathbf{F}}{\partial z_{ij}} \right\|_{\mathbb{F}}^2 + 2D^{-4} \frac{1}{n} \|\mathbf{F}\|_{\mathbb{F}}^2 n^{-1} D^2 [4 \max(1, (2\tau')^{-1})]^2 \end{aligned}$$

$$\begin{aligned}
&\leq 2D^{-2} \frac{1}{n} \sum_{ij} \left\| \frac{\partial \mathbf{F}}{\partial z_{ij}} \right\|_{\mathbb{F}}^2 + 2n^{-1} [4 \max(1, (2\tau')^{-1})]^2 \\
&\leq 4 \max(1, (\tau')^{-1} (T \wedge \frac{p}{n})) + 2n^{-1} [4 \max(1, (2\tau')^{-1})]^2 \\
&:= f(\tau', T, n, p),
\end{aligned}$$

where the second inequality is by Corollary E.3.

(2) For $\tau = 0$, by Lemma E.2, in the event $U_1 \cap U_2$, we obtain the same upper bounds as in the first case (1) with τ' replaced by η . To see this,

$$\begin{aligned}
&\|(\mathbf{I}_T \otimes \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}}) (\mathbf{F}^\top \otimes \mathbf{I}_p)\|_{\mathbb{F}}^2 \\
&\leq \|(\mathbf{I}_T \otimes \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}})\|_{\text{op}}^2 \|(\mathbf{F}^\top \otimes \mathbf{I}_p)\|_{\mathbb{F}}^2 \\
&= \|(\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}})\|_{\text{op}} p \|\mathbf{F}\|_{\mathbb{F}}^2 \\
&\leq \|(\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}})\|_{\text{op}} p \|\mathbf{F}\|_{\mathbb{F}}^2 \\
&\leq p \|\mathbf{F}\|_{\mathbb{F}}^2 \frac{1}{n\eta} \\
&= \frac{p}{n\eta} \|\mathbf{F}\|_{\mathbb{F}}^2,
\end{aligned}$$

where the third inequality is by Lemma B.2. Also, we have

$$\begin{aligned}
&\|(\mathbf{I}_T \otimes \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}}) (\mathbf{F}^\top \otimes \mathbf{I}_p)\|_{\mathbb{F}}^2 \\
&\leq \|(\mathbf{I}_T \otimes \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}})\|_{\mathbb{F}}^2 \|(\mathbf{F}^\top \otimes \mathbf{I}_p)\|_{\text{op}}^2 \\
&\leq \text{Tr} [(\mathbf{I}_T \otimes \mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \Sigma_{\hat{\mathcal{J}}, \hat{\mathcal{J}}}) \mathbf{M}^\dagger] \|\mathbf{F}\|_{\mathbb{F}}^2 \\
&\leq \text{Tr} [(\mathbf{I}_T \otimes \Sigma_{\hat{\mathcal{J}}, \hat{\mathcal{J}}}) \mathbf{M}^\dagger] \|\mathbf{F}\|_{\mathbb{F}}^2 \\
&\leq \text{Tr} [(\mathbf{I}_T \otimes \Sigma_{\hat{\mathcal{J}}, \hat{\mathcal{J}}}) (\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}})^\dagger)] \|\mathbf{F}\|_{\mathbb{F}}^2 \\
&= T \text{Tr} [\Sigma_{\hat{\mathcal{J}}, \hat{\mathcal{J}}} (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}})^\dagger] \|\mathbf{F}\|_{\mathbb{F}}^2 \\
&\leq T \text{Tr} [(n\eta)^{-1} \mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}})^\dagger] \|\mathbf{F}\|_{\mathbb{F}}^2 \\
&\leq \frac{T}{\eta} \|\mathbf{F}\|_{\mathbb{F}}^2,
\end{aligned}$$

where the penultimate inequality uses $\Sigma_{\hat{\mathcal{J}}, \hat{\mathcal{J}}} \preceq (n\eta)^{-1} \mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}}$ in the event $U_1 \cap U_2$. Therefore, on $U_1 \cap U_2$, we have

$$\begin{aligned}
\frac{1}{n} \sum_{ij} \left\| \frac{\partial \mathbf{F}/D}{\partial z_{ij}} \right\|_{\mathbb{F}}^2 &\leq 4 \max(1, (\eta)^{-1} (T \wedge \frac{p}{n})) + 2n^{-1} [4 \max(1, (2\tau')^{-1})]^2 \\
&:= f(\eta, T, n, p),
\end{aligned}$$

where the function f is the same as in case (1). The only difference is that τ' in the upper bound for case (1) is replaced by η in case (2).

■

G.2. Proofs of results in Section E.2

The following proof of Lemma E.6 relies on a similar argument as proof of Lemma E.4, we present the proof here for completeness.

Proof of Lemma E.6. Recall the KKT conditions for $\widehat{\mathbf{B}}$ defined in (6):

1) For $k \in \hat{\mathcal{J}}$, i.e., $\widehat{\mathbf{B}}^\top \mathbf{e}_k \neq \mathbf{0}$,

$$\mathbf{e}_k^\top \mathbf{X}^\top [\mathbf{E} - \mathbf{X}(\widehat{\mathbf{B}} - \mathbf{B}^*)] - n\tau \mathbf{e}_k^\top \widehat{\mathbf{B}} = \frac{n\lambda}{\|\widehat{\mathbf{B}}^\top \mathbf{e}_k\|_2} \mathbf{e}_k^\top \widehat{\mathbf{B}} \in \mathbb{R}^{1 \times T}.$$

2) For $k \notin \hat{\mathcal{J}}$, i.e., $\widehat{\mathbf{B}}^\top \mathbf{e}_k = \mathbf{0}$,

$$\left\| \mathbf{e}_k^\top \mathbf{X}^\top [\mathbf{E} - \mathbf{X}(\widehat{\mathbf{B}} - \mathbf{B}^*)] - n\tau \mathbf{e}_k^\top \widehat{\mathbf{B}} \right\| < n\lambda.$$

Here the strict inequality is guaranteed by Proposition 2.3 of Bellec (2020).

Let $\ddot{\mathbf{B}} = \frac{\partial \widehat{\mathbf{B}}}{\partial E_{it'}}$, $\dot{\mathbf{E}} = \frac{\partial \mathbf{E}}{\partial E_{it'}}$. Differentiation of the above display for $k \in \hat{\mathcal{J}}$ w.r.t. $E_{it'}$ yields

$$\mathbf{e}_k^\top \mathbf{X}^\top (\dot{\mathbf{E}} - \mathbf{X}\ddot{\mathbf{B}}) - n\tau \mathbf{e}_k^\top \ddot{\mathbf{B}} = n \mathbf{e}_k^\top \ddot{\mathbf{B}} \mathbf{H}^{(k)}$$

with $\mathbf{H}^{(k)} = \lambda \|\widehat{\mathbf{B}}^\top \mathbf{e}_k\|_2^{-1} \left(\mathbf{I}_T - \widehat{\mathbf{B}}^\top \mathbf{e}_k \mathbf{e}_k^\top \widehat{\mathbf{B}} \|\widehat{\mathbf{B}}^\top \mathbf{e}_k\|_2^{-2} \right) \in \mathbb{R}^{T \times T}$. Rearranging and using $\dot{\mathbf{E}} = \mathbf{e}_i \mathbf{e}_{t'}^\top$,

$$\mathbf{e}_k^\top \mathbf{X}^\top \mathbf{e}_i \mathbf{e}_{t'}^\top = \mathbf{e}_k^\top [n \ddot{\mathbf{B}} \mathbf{H}^{(k)} + (\mathbf{X}^\top \mathbf{X} + n\tau \mathbf{I}_{p \times p}) \ddot{\mathbf{B}}].$$

Recall $\mathbf{P}_{\hat{\mathcal{J}}} = \sum_{k \in \hat{\mathcal{J}}} \mathbf{e}_k \mathbf{e}_k^\top \in \mathbb{R}^{p \times p}$. Multiplying by \mathbf{e}_k to the left and summing over $k \in \hat{\mathcal{J}}$, we obtain

$$\mathbf{P}_{\hat{\mathcal{J}}} \mathbf{X}^\top \mathbf{e}_i \mathbf{e}_{t'}^\top = n \sum_{k \in \hat{\mathcal{J}}} \mathbf{e}_k \mathbf{e}_k^\top \ddot{\mathbf{B}} \mathbf{H}^{(k)} + \mathbf{P}_{\hat{\mathcal{J}}} (\mathbf{X}^\top \mathbf{X} + n\tau \mathbf{I}_{p \times p}) \ddot{\mathbf{B}},$$

which reduces to the following by $\text{supp}(\ddot{\mathbf{B}}) \subseteq \hat{\mathcal{J}}$ and $\mathbf{X}\ddot{\mathbf{B}} = \mathbf{X}_{\hat{\mathcal{J}}} \ddot{\mathbf{B}}$,

$$\begin{aligned} \mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{e}_i \mathbf{e}_{t'}^\top &= n \sum_{k \in \hat{\mathcal{J}}} \mathbf{e}_k \mathbf{e}_k^\top \ddot{\mathbf{B}} \mathbf{H}^{(k)} + \mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} \ddot{\mathbf{B}} \mathbf{I}_T + n\tau \mathbf{P}_{\hat{\mathcal{J}}} \ddot{\mathbf{B}} \mathbf{I}_T \\ &= n \sum_{k \in \hat{\mathcal{J}}} \mathbf{e}_k \mathbf{e}_k^\top \ddot{\mathbf{B}} \mathbf{H}^{(k)} + (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}}) \ddot{\mathbf{B}} \mathbf{I}_T. \end{aligned}$$

Vectorizing the above yields

$$\begin{aligned} (\mathbf{e}_{t'} \otimes \mathbf{X}_{\hat{\mathcal{J}}}^\top) \text{vec}(\mathbf{e}_i) &= [n \sum_{k \in \hat{\mathcal{J}}} (\mathbf{H}^{(k)} \otimes \mathbf{e}_k \mathbf{e}_k^\top) + \mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}})] \text{vec}(\ddot{\mathbf{B}}) \\ &= (n\tilde{\mathbf{H}} + \mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + n\tau \mathbf{P}_{\hat{\mathcal{J}}})) \text{vec}(\ddot{\mathbf{B}}) \\ &= \mathbf{M} \text{vec}(\ddot{\mathbf{B}}). \end{aligned}$$

A similar argument as in Proof of Lemma E.4 leads to

$$\text{vec}(\ddot{\mathbf{B}}) = \mathbf{M}^\dagger \mathbf{M} \text{vec}(\ddot{\mathbf{B}}) = \mathbf{M}^\dagger (\mathbf{e}_{t'} \otimes \mathbf{X}_{\hat{\mathcal{J}}}^\top) \mathbf{e}_i.$$

Therefore, by $\mathbf{X}\ddot{\mathbf{B}} = \mathbf{X}_{\hat{\mathcal{J}}}\ddot{\mathbf{B}}$,

$$\begin{aligned} \frac{\partial F_{lt}}{\partial E_{it'}} &= \mathbf{e}_l^\top \frac{\partial \mathbf{E} - \mathbf{X}(\hat{\mathbf{B}} - \mathbf{B}^*)}{\partial E_{it'}} \mathbf{e}_t \\ &= \mathbf{e}_l^\top (\mathbf{e}_i \mathbf{e}_{t'}^\top - \mathbf{X}\ddot{\mathbf{B}}) \mathbf{e}_t \\ &= \mathbf{e}_l^\top \mathbf{e}_i \mathbf{e}_{t'}^\top \mathbf{e}_t - \mathbf{e}_l^\top \mathbf{X}\ddot{\mathbf{B}} \mathbf{e}_t \\ &= \mathbf{e}_l^\top \mathbf{e}_i \mathbf{e}_{t'}^\top \mathbf{e}_t - (\mathbf{e}_t^\top \otimes \mathbf{e}_l^\top \mathbf{X}_{\hat{\mathcal{J}}}) \text{vec}(\ddot{\mathbf{B}}) \\ &= \mathbf{e}_l^\top \mathbf{e}_i \mathbf{e}_{t'}^\top \mathbf{e}_t - (\mathbf{e}_t^\top \otimes \mathbf{e}_l^\top \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger (\mathbf{e}_{t'} \otimes \mathbf{X}_{\hat{\mathcal{J}}}^\top) \mathbf{e}_i \\ &= \mathbf{e}_l^\top \mathbf{e}_i \mathbf{e}_{t'}^\top \mathbf{e}_t - \mathbf{e}_l^\top (\mathbf{e}_t^\top \otimes \mathbf{X}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger (\mathbf{e}_{t'} \otimes \mathbf{X}_{\hat{\mathcal{J}}}^\top) \mathbf{e}_i \\ &= \mathbf{e}_l^\top \mathbf{e}_i \mathbf{e}_{t'}^\top \mathbf{e}_t - \mathbf{e}_l^\top (\mathbf{e}_t^\top \otimes \mathbf{X}) \mathbf{M}^\dagger (\mathbf{e}_{t'} \otimes \mathbf{X}^\top) \mathbf{e}_i, \end{aligned}$$

where the last equality is due to $\mathbf{M}^\dagger = (\mathbf{I}_T \otimes \mathbf{P}_{\hat{\mathcal{J}}}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \mathbf{P}_{\hat{\mathcal{J}}})$.

Now the calculation of $\sum_{i=1}^n \frac{\partial F_{it}}{\partial E_{it'}}$ is straightforward,

$$\begin{aligned} \sum_{i=1}^n \frac{\partial F_{it}}{\partial E_{it'}} &= \sum_{i=1}^n [\mathbf{e}_i^\top \mathbf{e}_i \mathbf{e}_t^\top \mathbf{e}_{t'} - \mathbf{e}_i^\top (\mathbf{e}_t^\top \otimes \mathbf{X}) \mathbf{M}^\dagger (\mathbf{e}_{t'} \otimes \mathbf{X}^\top) \mathbf{e}_i] \\ &= n \mathbf{e}_t^\top \mathbf{e}_{t'} - \text{Tr}[(\mathbf{e}_t^\top \otimes \mathbf{X}) \mathbf{M}^\dagger (\mathbf{e}_{t'} \otimes \mathbf{X}^\top)] \\ &= n \mathbf{e}_t^\top \mathbf{e}_{t'} - \mathbf{e}_t^\top \hat{\mathbf{A}} \mathbf{e}_{t'} \\ &= \mathbf{e}_t^\top (n \mathbf{I}_T - \hat{\mathbf{A}}) \mathbf{e}_{t'}, \end{aligned}$$

where the third equality is due to the formula of $\hat{\mathbf{A}}$ in (10).

Noting that $\mathbf{F} = \mathbf{E} - \mathbf{Z}\mathbf{H}$, it follows that $\sum_{i=1}^n \frac{\partial \mathbf{e}_i^\top \mathbf{Z} \mathbf{H} \mathbf{e}_t}{\partial E_{it'}} = \mathbf{e}_t^\top \hat{\mathbf{A}} \mathbf{e}_{t'}$. ■

G.3. Proofs of results in Section E.3

Proof of Lemma E.9. Let $\mathbf{z} = \text{vec}(\mathbf{E})$, then $\mathbf{z} \sim \mathcal{N}(0, \mathbf{K})$ with $\mathbf{K} = \mathbf{S} \otimes \mathbf{I}_n$ by Assumption 1. For each $t_0, t'_0 \in [T]$, let $\mathbf{G}^{(t_0, t'_0)} = \mathbf{F} \mathbf{e}_{t'_0} \mathbf{e}_{t_0}^\top$, and $\mathbf{f}(\mathbf{z})^{(t_0, t'_0)} = \text{vec}(\mathbf{G}) \tilde{\mathbf{D}}^{-1}$. For convenience, we will drop the superscript (t_0, t'_0) from $\mathbf{G}^{(t_0, t'_0)}$ and $\mathbf{f}(\mathbf{z})^{(t_0, t'_0)}$ in this proof. By $\text{Tr}(\mathbf{A}^\top \mathbf{B}) = \text{vec}(\mathbf{A})^\top \text{vec}(\mathbf{B})$, we obtain

$$\mathbf{e}_{t_0}^\top \mathbf{E}^\top \mathbf{F} \tilde{\mathbf{D}}^{-1} \mathbf{e}_{t'_0} = \text{Tr}(\mathbf{E}^\top \mathbf{F} \mathbf{e}_{t'_0} \mathbf{e}_{t_0}^\top) \tilde{\mathbf{D}}^{-1} = \text{Tr}(\mathbf{E}^\top \mathbf{G} \tilde{\mathbf{D}}^{-1}) = \mathbf{z}^\top \mathbf{f}(\mathbf{z}). \quad (55)$$

By product rule, we have

$$\nabla \mathbf{f}(\mathbf{z}) = \frac{\partial \text{vec}(\mathbf{G})}{\partial \text{vec}(\mathbf{E})} \tilde{\mathbf{D}}^{-1} + \underbrace{\text{vec}(\mathbf{G}) \frac{\partial \tilde{\mathbf{D}}^{-1}}{\partial \text{vec}(\mathbf{E})}}_{\text{Rem}}, \quad (56)$$

where $\text{Rem} = \mathbf{u}\mathbf{v}^\top$ with $\mathbf{u} = \text{vec}(\mathbf{G}) \in \mathbb{R}^{nT \times 1}$, $\mathbf{v}^\top = \frac{\partial \tilde{D}^{-1}}{\partial \text{vec}(\mathbf{E})} \in \mathbb{R}^{1 \times nT}$. It follows that

$$\text{Tr}(\mathbf{K} \nabla f(\mathbf{z})) = \text{Tr} \left(\mathbf{K} \frac{\partial \text{vec}(\mathbf{G})}{\partial \text{vec}(\mathbf{E})} \right) \tilde{D}^{-1} + \text{Tr}(\mathbf{K} \text{Rem}). \quad (57)$$

Since $\mathbf{K} = \mathbf{S} \otimes \mathbf{I}_n$ and $\mathbf{G} = \mathbf{F} \mathbf{e}_{t'_0} \mathbf{e}_{t_0}^\top$, $\mathbf{K}_{it,lt'} = S_{it'} I(i=l)$, and $G_{it} = F_{it'} I(t=t_0)$. It follows

$$\text{Tr} \left(\mathbf{K} \frac{\partial \text{vec}(\mathbf{G})}{\partial \text{vec}(\mathbf{E})} \right) = \sum_{i,t} \sum_{l,t'} \mathbf{K}_{it,lt'} \frac{\partial G_{it}}{\partial E_{lt'}} = \sum_{t'} S_{t_0 t'} \sum_i \frac{\partial F_{it'_0}}{\partial E_{it'}} = \mathbf{e}_{t_0}^\top \mathbf{S} (n\mathbf{I}_T - \hat{\mathbf{A}}) \mathbf{e}_{t'_0}, \quad (58)$$

where the last equality used Lemma E.6 and that $\hat{\mathbf{A}}$ is symmetric.

Now we rewrite the quantity we want to bound as

$$\begin{aligned} & \mathbb{E} \left[\left\| \mathbf{E}^\top \mathbf{F} / \tilde{D} - \mathbf{S} (n\mathbf{I}_T - \hat{\mathbf{A}}) / \tilde{D} \right\|_{\mathbb{F}}^2 \right] \\ &= \sum_{t_0, t'_0} \mathbb{E} \left[\left(\mathbf{e}_{t_0}^\top \mathbf{E}^\top \mathbf{F} \tilde{D}^{-1} \mathbf{e}_{t'_0} - \mathbf{e}_{t_0}^\top \mathbf{S} (n\mathbf{I}_T - \hat{\mathbf{A}}) \mathbf{e}_{t'_0} \tilde{D}^{-1} \right)^2 \right] \\ &= \sum_{t_0, t'_0} \mathbb{E} \left[\left(\mathbf{z}^\top f(\mathbf{z}) - \text{Tr}(\mathbf{K} \nabla f(\mathbf{z})) + \text{Tr}(\mathbf{K} \text{Rem}) \right)^2 \right] \\ &\leq 2 \sum_{t_0, t'_0} \left\{ \mathbb{E} \left[\left(\mathbf{z}^\top f(\mathbf{z}) - \text{Tr}(\mathbf{K} \nabla f(\mathbf{z})) \right)^2 \right] + \mathbb{E} \left[\left(\text{Tr}(\mathbf{K} \text{Rem}) \right)^2 \right] \right\}, \end{aligned} \quad (59)$$

where the second equality follows from (55), (57) and (58), and the last inequality uses elementary inequality $(a+b)^2 \leq 2(a^2+b^2)$. We next bound the two terms in (59). **First term in (59).** By second-order Stein formula in Lemma E.7,

$$\sum_{t_0, t'_0} \mathbb{E} \left(\mathbf{z}^\top f(\mathbf{z}) - \text{Tr}(\mathbf{K} \nabla f(\mathbf{z})) \right)^2 = \sum_{t_0, t'_0} \mathbb{E} \left[\left\| \mathbf{K}^{\frac{1}{2}} f(\mathbf{z}) \right\|_{\mathbb{F}}^2 + \text{Tr} \left[(\mathbf{K} \nabla f(\mathbf{z}))^2 \right] \right]. \quad (60)$$

Now we bound the two terms in the right-hand side of (60). For the first term, recall $f(\mathbf{z}) = \text{vec}(\mathbf{G}) \tilde{D}^{-1}$, and $\mathbf{G} = \mathbf{F} \mathbf{e}_{t'_0} \mathbf{e}_{t_0}^\top$, we obtain

$$\left\| \mathbf{K}^{\frac{1}{2}} f(\mathbf{z}) \right\|_{\mathbb{F}}^2 = \tilde{D}^{-2} \left\| (\mathbf{S}^{\frac{1}{2}} \otimes \mathbf{I}_n) \text{vec}(\mathbf{G}) \right\|_{\mathbb{F}}^2 = \tilde{D}^{-2} \left\| \mathbf{G} \mathbf{S}^{\frac{1}{2}} \right\|_{\mathbb{F}}^2 = \tilde{D}^{-2} \left\| \mathbf{S}^{\frac{1}{2}} \mathbf{e}_{t_0} \right\|_{\mathbb{F}}^2 \left\| \mathbf{F} \mathbf{e}_{t'_0} \right\|_{\mathbb{F}}^2.$$

Summing over all $(t_0, t'_0) \in [T] \times [T]$, we obtain

$$\sum_{t_0, t'_0} \left\| \mathbf{K}^{\frac{1}{2}} f(\mathbf{z}) \right\|_{\mathbb{F}}^2 = \tilde{D}^{-2} \left\| \mathbf{F} \right\|_{\mathbb{F}}^2 \text{Tr}(\mathbf{S}). \quad (61)$$

For the second term in RHS of (60), recall $\nabla f(\mathbf{z}) = \frac{\partial \text{vec}(\mathbf{G})}{\partial \text{vec}(\mathbf{E})} \tilde{D}^{-1} + \text{Rem}$,

$$\begin{aligned} & \text{Tr} \left[(\mathbf{K} \nabla f(\mathbf{z}))^2 \right] \\ &= \tilde{D}^{-2} \text{Tr} \left[\left(\mathbf{K} \frac{\partial \text{vec}(\mathbf{G})}{\partial \text{vec}(\mathbf{E})} \right)^2 \right] + \text{Tr}[(\mathbf{K} \text{Rem})^2] + 2\tilde{D}^{-1} \text{Tr} \left[\mathbf{K} \frac{\partial \text{vec}(\mathbf{G})}{\partial \text{vec}(\mathbf{E})} \mathbf{K} \text{Rem} \right]. \end{aligned} \quad (62)$$

By property of vectorization operation, $\mathbf{vec}(G) = \mathbf{vec}(F e_{t'_0} e_{t_0}^\top) = (e_{t_0} e_{t'_0}^\top \otimes I_n) \mathbf{vec}(F)$, hence

$$\frac{\partial \mathbf{vec}(G)}{\partial \mathbf{vec}(E)} = (e_{t_0} e_{t'_0}^\top \otimes I_n) \frac{\partial \mathbf{vec}(F)}{\partial \mathbf{vec}(E)},$$

where $\|\frac{\partial \mathbf{vec}(F)}{\partial \mathbf{vec}(E)}\|_{\text{op}} \leq 1$ since the map $\mathbf{vec}(E) \mapsto \mathbf{vec}(F)$ is 1-Lipschitz by (Bellec and Tsybakov, 2016, proposition 3).

Now we bound the three terms in (62). For the first term, by Cauchy-Schwarz inequality,

$$\begin{aligned} & \tilde{D}^{-2} \text{Tr} \left[\left(K \frac{\partial \mathbf{vec}(G)}{\partial \mathbf{vec}(E)} \right)^2 \right] \\ &= \tilde{D}^{-2} \text{Tr} \left(K (e_{t_0} e_{t'_0}^\top \otimes I_n) \frac{\partial \mathbf{vec}(F)}{\partial \mathbf{vec}(E)} K (e_{t_0} e_{t'_0}^\top \otimes I_n) \frac{\partial \mathbf{vec}(F)}{\partial \mathbf{vec}(E)} \right) \\ &\leq \tilde{D}^{-2} \|(e_{t'_0}^\top \otimes I_n) \frac{\partial \mathbf{vec}(F)}{\partial \mathbf{vec}(E)} K (e_{t_0} \otimes I_n)\|_{\text{F}}^2. \end{aligned}$$

For the second term in (62), recall $\text{Rem} = \mathbf{u} \mathbf{v}^\top$, and $\mathbf{u} = \mathbf{vec}(G)$, $\mathbf{v}^\top = \frac{\partial \tilde{D}^{-1}}{\partial \mathbf{vec}(E)}$ from (56), then $\text{Tr}[(K \text{Rem})^2] = \text{Tr}(K \mathbf{u} \mathbf{v}^\top K \mathbf{u} \mathbf{v}^\top) = (\mathbf{v}^\top K \mathbf{u})^2$, thus,

$$\begin{aligned} & \text{Tr}[(K \text{Rem})^2] \\ &= \left[\frac{\partial \tilde{D}^{-1}}{\partial \mathbf{vec}(E)} K \mathbf{vec}(G) \right]^2 \\ &= \tilde{D}^{-6} \left[\mathbf{vec}(F)^\top \frac{\partial \mathbf{vec}(F)}{\partial \mathbf{vec}(E)} K (e_{t_0} e_{t'_0}^\top \otimes I_n) \mathbf{vec}(F) \right]^2 \\ &\leq \tilde{D}^{-6} \|\mathbf{vec}(F)^\top \frac{\partial \mathbf{vec}(F)}{\partial \mathbf{vec}(E)} K (e_{t_0} \otimes I_n)\|^2 \|(e_{t'_0}^\top \otimes I_n) \mathbf{vec}(F)\|^2. \end{aligned}$$

where the second equality uses $\frac{\partial \tilde{D}^{-1}}{\partial \mathbf{vec}(E)} = -\frac{1}{2} \tilde{D}^{-3} \frac{\partial \|\mathbf{F}\|_{\text{F}}^2}{\partial \mathbf{vec}(E)} = -\tilde{D}^{-3} \mathbf{vec}(F)^\top \frac{\partial \mathbf{vec}(F)}{\partial \mathbf{vec}(E)}$ by chain rule, and the inequality uses Cauchy-Schwarz inequality.

For the third term in (62), recall $\text{Rem} = \mathbf{u} \mathbf{v}^\top$, and $\mathbf{u} = \mathbf{vec}(G)$, $\mathbf{v}^\top = \frac{\partial \tilde{D}^{-1}}{\partial \mathbf{vec}(E)}$, then we have $\text{Tr} \left[K \frac{\partial \mathbf{vec}(G)}{\partial \mathbf{vec}(E)} K \text{Rem} \right] = \mathbf{v}^\top K \frac{\partial \mathbf{vec}(G)}{\partial \mathbf{vec}(E)} K \mathbf{u}$, hence

$$\begin{aligned} & 2\tilde{D}^{-1} \text{Tr} \left[K \frac{\partial \mathbf{vec}(G)}{\partial \mathbf{vec}(E)} K \text{Rem} \right] \\ &= 2\tilde{D}^{-1} \frac{\partial \tilde{D}^{-1}}{\partial \mathbf{vec}(E)} K \frac{\partial \mathbf{vec}(G)}{\partial \mathbf{vec}(E)} K \mathbf{vec}(G) \\ &= -2\tilde{D}^{-4} \mathbf{vec}(F)^\top \frac{\partial \mathbf{vec}(F)}{\partial \mathbf{vec}(E)} K (e_{t_0} e_{t'_0}^\top \otimes I_n) \frac{\partial \mathbf{vec}(F)}{\partial \mathbf{vec}(E)} K (e_{t_0} e_{t'_0}^\top \otimes I_n) \mathbf{vec}(F) \\ &\leq 2\tilde{D}^{-4} \|(e_{t'_0}^\top \otimes I_n) \mathbf{vec}(F) \mathbf{vec}(F)^\top \frac{\partial \mathbf{vec}(F)}{\partial \mathbf{vec}(E)} K (e_{t_0} \otimes I_n)\|_{\text{F}} \\ &\quad \|(e_{t'_0}^\top \otimes I_n) \frac{\partial \mathbf{vec}(F)}{\partial \mathbf{vec}(E)} K (e_{t_0} \otimes I_n)\|_{\text{F}}, \end{aligned}$$

where the last inequality uses Cauchy-Schwarz inequality.

Summing over all $(t_0, t'_0) \in [T] \times [T]$ for these three terms in (62), using $\|\frac{\partial \text{vec}(\mathbf{F})}{\partial \text{vec}(\mathbf{E})}\|_{\text{op}} \leq 1$, $\mathbf{K} = \mathbf{S} \otimes \mathbf{I}_n$, and $\|\mathbf{S}\|_{\text{F}} \leq \|\mathbf{S}^{\frac{1}{2}}\|_{\text{F}}^2 = \text{Tr}(\mathbf{S})$, we obtain

$$\begin{aligned}
& \sum_{t_0, t'_0} \text{Tr}[(\mathbf{K} \nabla f(\mathbf{z}))^2] \\
& \leq \tilde{D}^{-2} \|\mathbf{K}\|_{\text{F}}^2 + \tilde{D}^{-6} \|\mathbf{F}\|_{\text{F}}^2 \|\mathbf{K}\|_{\text{F}}^2 \|\mathbf{F}\|_{\text{F}}^2 + 2\tilde{D}^{-4} \|\mathbf{F}\|_{\text{F}}^2 \|\mathbf{K}\|_{\text{F}}^2 \\
& = \tilde{D}^{-2} n \|\mathbf{S}\|_{\text{F}}^2 + \tilde{D}^{-6} \|\mathbf{F}\|_{\text{F}}^4 n \|\mathbf{S}\|_{\text{F}}^2 + 2\tilde{D}^{-4} \|\mathbf{F}\|_{\text{F}}^2 n \|\mathbf{S}\|_{\text{F}}^2 \\
& \leq [\tilde{D}^{-2} n \text{Tr}(\mathbf{S}) + \tilde{D}^{-6} \|\mathbf{F}\|_{\text{F}}^4 n \text{Tr}(\mathbf{S}) + 2\tilde{D}^{-4} \|\mathbf{F}\|_{\text{F}}^2 n \text{Tr}(\mathbf{S})] \text{Tr}(\mathbf{S}). \tag{63}
\end{aligned}$$

Second term in (59). Recall that $\text{Rem} = \mathbf{u}\mathbf{v}^\top$, hence $[\text{Tr}(\mathbf{K}\text{Rem})]^2 = \text{Tr}[(\mathbf{K}\text{Rem})^2]$. By calculation of second term in (62), we obtain

$$\sum_{t_0, t'_0} [\text{Tr}(\mathbf{K}\text{Rem})]^2 = \sum_{t_0, t'_0} \text{Tr}[(\mathbf{K}\text{Rem})^2] \leq \tilde{D}^{-6} \|\mathbf{F}\|_{\text{F}}^4 n \text{Tr}(\mathbf{S})^2. \tag{64}$$

Combining the results (59), (60), (61), (63), (64), we obtain

$$\begin{aligned}
& \mathbb{E} \left[\left\| \mathbf{E}^\top \mathbf{F} / \tilde{D} - \mathbf{S}(n\mathbf{I}_T - \hat{\mathbf{A}}) / \tilde{D} \right\|_{\text{F}}^2 \right] \\
& \leq 2[\tilde{D}^{-2} \|\mathbf{F}\|_{\text{F}}^2 + \tilde{D}^{-2} n \text{Tr}(\mathbf{S}) + 2\tilde{D}^{-6} \|\mathbf{F}\|_{\text{F}}^4 n \text{Tr}(\mathbf{S}) + 2\tilde{D}^{-4} \|\mathbf{F}\|_{\text{F}}^2 n \text{Tr}(\mathbf{S})] \text{Tr}(\mathbf{S}) \\
& \leq 4 \text{Tr}(\mathbf{S}),
\end{aligned}$$

thanks to $\tilde{D}^2 = \|\mathbf{F}\|_{\text{F}}^2 + n \text{Tr}(\mathbf{S})$. ■

Proof of Lemma E.10. Apply (Bellec, 2020, Proposition 6.3) with $\boldsymbol{\rho} = \mathbf{U}\mathbf{e}_t$, $\boldsymbol{\eta} = \mathbf{V}\mathbf{e}_{t'}$, we obtain

$$\begin{aligned}
& \mathbb{E} \left[\left\| \mathbf{U}^\top \mathbf{Z}\mathbf{V} - \sum_{j=1}^p \sum_{i=1}^n \frac{\partial}{\partial z_{ij}} (\mathbf{U}^\top \mathbf{e}_i \mathbf{e}_j^\top \mathbf{V}) \right\|_{\text{F}}^2 \right] \\
& = \sum_{t, t'=1}^T \mathbb{E} (\mathbf{e}_t^\top \mathbf{U}^\top \mathbf{Z}\mathbf{V} \mathbf{e}_{t'} - \sum_{j=1}^p \sum_{i=1}^n \frac{\partial}{\partial z_{ij}} \mathbf{e}_t^\top \mathbf{U}^\top \mathbf{e}_i \mathbf{e}_j^\top \mathbf{V} \mathbf{e}_{t'})^2 \\
& \leq \sum_{t, t'=1}^T \left[\mathbb{E} [\|\mathbf{U}\mathbf{e}_t\|^2 \|\mathbf{V}\mathbf{e}_{t'}\|^2] + \mathbb{E} \sum_{ij} \left[2\|\mathbf{V}\mathbf{e}_{t'}\|^2 \left\| \frac{\partial \mathbf{U}\mathbf{e}_t}{\partial z_{ij}} \right\|_{\text{F}} + 2\|\mathbf{U}\mathbf{e}_t\|^2 \left\| \frac{\partial \mathbf{V}\mathbf{e}_{t'}}{\partial z_{ij}} \right\|_{\text{F}} \right] \right] \\
& = \mathbb{E} \|\mathbf{U}\|_{\text{F}}^2 \|\mathbf{V}\|_{\text{F}}^2 + \mathbb{E} \sum_{ij} \left[2\|\mathbf{V}\|_{\text{F}}^2 \left\| \frac{\partial \mathbf{U}}{\partial z_{ij}} \right\|_{\text{F}}^2 + 2\|\mathbf{U}\|_{\text{F}}^2 \left\| \frac{\partial \mathbf{V}}{\partial z_{ij}} \right\|_{\text{F}}^2 \right].
\end{aligned}$$

■

Proof of Corollary E.11. By Kirszbraun's theorem, there exists an L_1 -Lipschitz function $\tilde{\mathbf{U}} : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^{n \times T}$ such that $\tilde{\mathbf{U}} = \mathbf{U}$ on Ω , and an L_2 -Lipschitz function $\tilde{\mathbf{V}} : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^{n \times T}$ such that $\tilde{\mathbf{V}} = \mathbf{V}$ on Ω . By

projecting \bar{U} and \bar{V} onto the Euclidean ball of radius 1 and K if necessary, we assume without loss of generality that $\|\bar{U}\|_F \leq 1$ and $\|\bar{V}\|_F \leq K$. Therefore,

$$\begin{aligned}
& \mathbb{E} \left[I(\Omega) \left\| \mathbf{U}^\top \mathbf{Z} \mathbf{V} - \sum_{j=1}^p \sum_{i=1}^n \frac{\partial}{\partial z_{ij}} (\mathbf{U}^\top \mathbf{e}_i \mathbf{e}_j^\top \mathbf{V}) \right\|_F^2 \right] \\
&= \mathbb{E} \left[I(\Omega) \left\| \bar{\mathbf{U}}^\top \mathbf{Z} \bar{\mathbf{V}} - \sum_{j=1}^p \sum_{i=1}^n \frac{\partial}{\partial z_{ij}} (\bar{\mathbf{U}}^\top \mathbf{e}_i \mathbf{e}_j^\top \bar{\mathbf{V}}) \right\|_F^2 \right] \\
&\leq \mathbb{E} \left[\left\| \bar{\mathbf{U}}^\top \mathbf{Z} \bar{\mathbf{V}} - \sum_{j=1}^p \sum_{i=1}^n \frac{\partial}{\partial z_{ij}} (\bar{\mathbf{U}}^\top \mathbf{e}_i \mathbf{e}_j^\top \bar{\mathbf{V}}) \right\|_F^2 \right] \\
&\leq \mathbb{E} \left[(\|\bar{\mathbf{U}}\|_F^2 \|\bar{\mathbf{V}}\|_F^2) + 2 \sum_{ij} \left(\|\bar{\mathbf{V}}\|_F^2 \left\| \frac{\partial \bar{\mathbf{U}}}{\partial z_{ij}} \right\|_F^2 + \|\bar{\mathbf{U}}\|_F^2 \left\| \frac{\partial \bar{\mathbf{V}}}{\partial z_{ij}} \right\|_F^2 \right) \right] \\
&\leq K^2 + 2\mathbb{E} \left[\sum_{ij} \left(K^2 \left\| \frac{\partial \bar{\mathbf{U}}}{\partial z_{ij}} \right\|_F^2 + \left\| \frac{\partial \bar{\mathbf{V}}}{\partial z_{ij}} \right\|_F^2 \right) \right] \\
&= K^2 + 2\mathbb{E} \left[I(\Omega) \sum_{ij} \left(K^2 \left\| \frac{\partial \bar{\mathbf{U}}}{\partial z_{ij}} \right\|_F^2 + \left\| \frac{\partial \bar{\mathbf{V}}}{\partial z_{ij}} \right\|_F^2 \right) + I(\Omega^c) \sum_{ij} \left(K^2 \left\| \frac{\partial \bar{\mathbf{U}}}{\partial z_{ij}} \right\|_F^2 + \left\| \frac{\partial \bar{\mathbf{V}}}{\partial z_{ij}} \right\|_F^2 \right) \right] \\
&= K^2 + 2\mathbb{E} \left[I(\Omega) \sum_{ij} \left(K^2 \left\| \frac{\partial \mathbf{U}}{\partial z_{ij}} \right\|_F^2 + \left\| \frac{\partial \mathbf{V}}{\partial z_{ij}} \right\|_F^2 \right) + I(\Omega^c) \sum_{ij} \left(K^2 \left\| \frac{\partial \bar{\mathbf{U}}}{\partial z_{ij}} \right\|_F^2 + \left\| \frac{\partial \bar{\mathbf{V}}}{\partial z_{ij}} \right\|_F^2 \right) \right] \\
&\leq K^2 + 2\mathbb{E} \left[I(\Omega) \sum_{ij} \left(K^2 \left\| \frac{\partial \mathbf{U}}{\partial z_{ij}} \right\|_F^2 + \left\| \frac{\partial \mathbf{V}}{\partial z_{ij}} \right\|_F^2 \right) \right] + 2C(K^2 L_1^2 + L_2^2),
\end{aligned}$$

where the last inequality uses $\sum_{ij} \left\| \frac{\partial \bar{\mathbf{U}}}{\partial z_{ij}} \right\|_F^2 \leq nT(n^{-1/2}L_1)^2 = TL_1^2$, $\sum_{ij} \left\| \frac{\partial \bar{\mathbf{V}}}{\partial z_{ij}} \right\|_F^2 \leq TL_2^2$ by Lipschitz properties of $\bar{\mathbf{U}}$, $\bar{\mathbf{V}}$, and $P(\Omega^c) \leq C/T$. ■

Proof of Lemma E.12. For each $j \in [p]$, let $\mathbb{E}_j(\cdot)$ denote the conditional expectation $\mathbb{E}[\cdot | \{\mathbf{Z} \mathbf{e}_k, k \neq j\}]$. The left-hand side of the desired inequality can be rewritten as

$$\mathbb{E} \left[\left\| p \mathbf{U}^\top \mathbf{V} - \sum_{j=1}^p (\mathbb{E}_j \mathbf{U}^\top \mathbf{Z} - \mathbf{L}^\top) \mathbf{e}_j \mathbf{e}_j^\top (\mathbf{Z}^\top \mathbb{E}_j \mathbf{V} - \hat{\mathbf{L}}) \right\|_F \right]$$

with $\mathbf{L} \in \mathbb{R}^{p \times T}$ defined by $\mathbf{L}^\top \mathbf{e}_j = \mathbb{E}_j \mathbf{U}^\top \mathbf{Z} \mathbf{e}_j - \mathbf{U}^\top \mathbf{Z} \mathbf{e}_j + \sum_{i=1}^n \partial_{ij} \mathbf{U}^\top \mathbf{e}_i$ and $\hat{\mathbf{L}}$ defined similarly with \mathbf{U} replaced by \mathbf{V} . We develop the terms in the sum over j as follows:

$$\begin{aligned}
& p \mathbf{U}^\top \mathbf{V} - \sum_j (\mathbb{E}_j \mathbf{U}^\top \mathbf{Z} - \mathbf{L}^\top) \mathbf{e}_j \mathbf{e}_j^\top (\mathbf{Z}^\top \mathbb{E}_j \mathbf{V} - \hat{\mathbf{L}}) \\
&= \sum_j \left(\mathbf{U}^\top \mathbf{V} - \mathbb{E}_j [\mathbf{U}^\top] \mathbb{E}_j \mathbf{V} \right) \tag{65}
\end{aligned}$$

$$+ \sum_j \left(\mathbb{E}_j [\mathbf{U}^\top] \mathbb{E}_j \mathbf{V} - \mathbb{E}_j \mathbf{U}^\top \mathbf{Z} \mathbf{e}_j \mathbf{e}_j^\top \mathbf{Z}^\top \mathbb{E}_j \mathbf{V} \right) \tag{66}$$

$$- \mathbf{L}^\top \hat{\mathbf{L}} \quad (67)$$

$$+ \sum_j \left(\mathbb{E}_j \mathbf{U}^\top \mathbf{Z} \mathbf{e}_j \mathbf{e}_j^\top \hat{\mathbf{L}} \right) + \left(\mathbf{L}^\top \mathbf{e}_j \mathbf{e}_j^\top \mathbf{Z}^\top \mathbb{E}_j[\mathbf{U}] \right). \quad (68)$$

First, for (67), by the Cauchy-Schwarz inequality $\mathbb{E}[\|\mathbf{L}^\top \hat{\mathbf{L}}\|_F] \leq \mathbb{E}[\|\mathbf{L}\|_F^2]^{\frac{1}{2}} \mathbb{E}[\|\hat{\mathbf{L}}\|_F^2]^{\frac{1}{2}}$. For a fixed $j \in [p]$ and $t \in [T]$,

$$\begin{aligned} \mathbb{E}[(\mathbf{e}_j^\top \mathbf{L} \mathbf{e}_t)^2] &\leq \sum_{i=1}^n \mathbb{E}[(\mathbf{e}_i^\top (\mathbb{E}_j[\mathbf{U}] - \mathbf{U}) \mathbf{e}_t)^2] + \mathbb{E} \sum_{i=1}^n \sum_{l=1}^n \left(\frac{\mathbf{e}_i^\top \partial \mathbf{U} \mathbf{e}_t}{\partial z_{lj}} \right)^2 \\ &\leq 2 \mathbb{E} \sum_{i=1}^n \sum_{l=1}^n \left(\frac{\mathbf{e}_i^\top \partial \mathbf{U} \mathbf{e}_t}{\partial z_{lj}} \right)^2, \end{aligned}$$

where the two inequalities are due to the second-order stein inequality in Lemma E.7, and Gaussian-Poincaré inequality in Lemma E.8, respectively. Summing over $j \in [p]$ and $t \in [T]$ we obtain $\mathbb{E}[\|\mathbf{L}\|_F^2] \leq 2 \mathbb{E} \sum_{l,j} \|\partial_{lj} \mathbf{U}\|_F^2 = 2 \|\mathbf{U}\|_\partial^2$. Combined with the same bound for $\hat{\mathbf{L}}$, we obtain $\mathbb{E}[\|(\text{67})\|_F^2] \leq 2 \|\mathbf{U}\|_\partial \|\mathbf{V}\|_\partial$. We now turn to the two terms in (68). By the triangle inequality for the Frobenius norm,

$$\begin{aligned} \mathbb{E} \left[\left\| \sum_j \mathbb{E}_j \mathbf{U}^\top \mathbf{Z} \mathbf{e}_j \mathbf{e}_j^\top \hat{\mathbf{L}} \right\|_F \right] &\leq \sum_j \mathbb{E} \left[\|\mathbb{E}_j \mathbf{U}^\top \mathbf{Z} \mathbf{e}_j\|_2 \|\mathbf{e}_j^\top \hat{\mathbf{L}}\|_2 \right] \\ &\leq \mathbb{E} \left[\sum_j \|\mathbb{E}_j \mathbf{U}^\top \mathbf{Z} \mathbf{e}_j\|_2^2 \right]^{\frac{1}{2}} \mathbb{E} \left[\sum_j \|\mathbf{e}_j^\top \hat{\mathbf{L}}\|_2^2 \right]^{\frac{1}{2}} \\ &\leq (p \mathbb{E}[\|\mathbf{U}\|_F^2])^{\frac{1}{2}} \mathbb{E}[\|\hat{\mathbf{L}}\|_F^2]^{\frac{1}{2}}, \end{aligned}$$

where we used that $\|\mathbf{a} \mathbf{b}^\top\|_F = \|\mathbf{a}\|_2 \|\mathbf{b}\|_2$ for two vectors \mathbf{a}, \mathbf{b} , the Cauchy-Schwarz inequality, $\mathbb{E}[\|\mathbf{A} \mathbf{z}_j\|_2^2 | \mathbf{A}] = \|\mathbf{A}\|_F^2$ if matrix \mathbf{A} is independent of $\mathbf{z}_j \sim \mathcal{N}(0, I_n)$ (set $\mathbf{z}_j = \mathbf{Z} \mathbf{e}_j$), and Jensen's inequality.

Next, we decompose (65) as $\sum_j \mathbf{U}^\top (\mathbf{V} - \mathbb{E}_j \mathbf{V}) + \sum_j (\mathbf{U} - \mathbb{E}_j \mathbf{U})^\top \mathbb{E}_j \mathbf{V}$. We have by the submultiplicativity of the Frobenius norm and the Cauchy-Schwarz inequality

$$\begin{aligned} \mathbb{E}[\|\mathbf{U}^\top \sum_j (\mathbf{V} - \mathbb{E}_j \mathbf{V})\|_F] &\leq \mathbb{E} \left[\sum_j \|\mathbf{U}\|_F \|\mathbf{V} - \mathbb{E}_j \mathbf{V}\|_F \right] \\ &\leq \mathbb{E}[p \|\mathbf{U}\|_F^2]^{\frac{1}{2}} \mathbb{E} \left[\sum_j \|\mathbf{V} - \mathbb{E}_j \mathbf{V}\|_F^2 \right]^{\frac{1}{2}}. \end{aligned}$$

By the Gaussian Poincaré inequality applied p times, $\mathbb{E}[\sum_j \|\mathbf{V} - \mathbb{E}_j \mathbf{V}\|_F^2] \leq \|\mathbf{V}\|_\partial^2$, so that the previous display is bounded from above by $\sqrt{p} \|\mathbf{V}\|_\partial$. Similarly, $\mathbb{E}[\|\sum_j (\mathbb{E}_j[\mathbf{U}] - \mathbf{U})^\top \mathbb{E}_j \mathbf{V}\|_F] \leq \sqrt{p} \|\mathbf{U}\|_\partial$ and $\mathbb{E}[\|(\text{65})\|_F] \leq \sqrt{p} (\|\mathbf{U}\|_\partial + \|\mathbf{V}\|_\partial)$.

For the last remaining term, (66), we first use $\mathbb{E}[\|(\text{66})\|_F] \leq \mathbb{E}[\|(\text{66})\|_F^2]^{\frac{1}{2}}$ by Jensen's inequality and now proceed to bound $\|(\text{66})\|_F^2$. We have

$$\|(\text{66})\|_F^2 = \left\| \sum_j \mathbb{E}_j \mathbf{U}^\top \mathbb{E}_j \mathbf{V} - \mathbb{E}_j \mathbf{U}^\top \mathbf{X} \mathbf{e}_j \mathbf{e}_j^\top \mathbf{Z}^\top \mathbb{E}_j \mathbf{V} \right\|_F^2 = \sum_{j,k} \text{Tr}[\mathbf{M}_j^\top \mathbf{M}_k],$$

where $\mathbf{M}_j = \mathbb{E}_j \mathbf{U}^\top \mathbb{E}_j \mathbf{V} - \mathbb{E}_j \mathbf{U}^\top \mathbf{Z} \mathbf{e}_j \mathbf{e}_j^\top \mathbf{Z}^\top \mathbb{E}_j \mathbf{V}$. We first bound $\sum_j \|\mathbf{M}_j\|_F^2$. Since the variance of $\mathbf{a}^\top \mathbf{b} - \mathbf{a}^\top \mathbf{g} \mathbf{g}^\top \mathbf{b}$ for standard normal $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_p)$ is $2\|(\mathbf{a}\mathbf{b}^\top + \mathbf{b}\mathbf{a}^\top)/2\|_F^2 \leq 2\|\mathbf{a}\|_2^2 \|\mathbf{b}\|_2^2$, applying this variance bound on each pair of coordinates $(t, t') \in [T] \times [T]$ gives $\sum_j \|\mathbf{M}_j\|_F^2 \leq \sum_j 2\|\mathbb{E}_j[\mathbf{U}]\|_F^2 \|\mathbb{E}_j[\mathbf{V}]\|_F^2 \leq 2p$.

We now bound $\sum_{j \neq k} \text{Tr}[\mathbf{M}_j^\top \mathbf{M}_k]$. Setting $\mathbf{z}_j = \mathbf{Z} \mathbf{e}_j \sim \mathcal{N}(0, \mathbf{I}_n)$ for every $j \in [p]$, we will use many times the identity

$$\mathbb{E}[(\mathbf{z}_j^\top f(\mathbf{Z}) - \sum_i \partial_{ij} f(\mathbf{Z})^\top \mathbf{e}_i) g(\mathbf{Z})] = \mathbb{E}[\sum_i f(\mathbf{Z})^\top \mathbf{e}_i \partial_{ij} g(\mathbf{Z})] \quad (69)$$

which follows from Stein's formula for $f : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^n$ and $g : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}$. With $f^{tt'}(\mathbf{Z}) = (\mathbf{z}_j^\top \mathbb{E}_j[\mathbf{U}] \mathbf{e}_{t'}) \mathbb{E}_j \mathbf{V} e_t$ and $g^{tt'}(\mathbf{Z}) = \mathbf{e}_{t'}^\top \mathbf{M}_k e_t$, we find

$$\begin{aligned} \mathbb{E} \text{Tr}[\mathbf{M}_j^\top \mathbf{M}_k] &= \mathbb{E} \text{Tr}[\mathbf{M}_j^\top \sum_t \mathbf{e}_{t'} \mathbf{e}_{t'}^\top \mathbf{M}_k] = \mathbb{E}[\sum_{tt'} \mathbf{e}_t^\top \mathbf{M}_j^\top \mathbf{e}_{t'} \mathbf{e}_{t'}^\top \mathbf{M}_k e_t] \\ &= \mathbb{E} \sum_{tt'} \left(\mathbf{z}_j^\top f^{tt'}(\mathbf{Z}) - \sum_i \mathbf{e}_i^\top \partial_{ij} f^{tt'}(\mathbf{Z}) \right) g^{tt'}(\mathbf{Z}) \\ &= \mathbb{E} \sum_{tt'} \sum_i \mathbf{e}_i^\top f^{tt'}(\mathbf{Z}) \partial_{ij} g^{tt'}(\mathbf{Z}). \end{aligned}$$

where $g_{tt'}(\mathbf{Z}) = (\mathbf{e}_t^\top \mathbb{E}_k \mathbf{V}^\top \mathbb{E}_k \mathbf{V} e_{t'} - \mathbf{e}_t^\top \mathbb{E}_k \mathbf{U}^\top \mathbf{z}_k \mathbf{z}_k^\top \mathbb{E}_k \mathbf{V} e_{t'})$ and

$$\partial_{ij} g_{tt'} = \partial_{ij} \mathbf{e}_{t'}^\top \mathbf{M}_k e_t = \mathbf{e}_{t'}^\top \partial_{ij} [\mathbb{E}_k \mathbf{U}^\top \mathbb{E}_k \mathbf{V}] e_t - \mathbf{z}_k^\top \partial_{ij} [\mathbb{E}_k \mathbf{U} \mathbf{e}_{t'} \mathbf{e}_{t'}^\top \mathbb{E}_k \mathbf{U}^\top] \mathbf{z}_k.$$

Now define $\tilde{f}^{tt'}(\mathbf{Z}) = \partial_{ij} [\mathbb{E}_k \mathbf{U} \mathbf{e}_{t'} \mathbf{e}_{t'}^\top \mathbb{E}_k \mathbf{U}^\top] \mathbf{z}_k$ and $\tilde{g}^{tt'}(\mathbf{Z}) = \sum_i \mathbf{e}_i^\top f^{tt'}(\mathbf{Z})$. By definition of $\tilde{f}^{tt'}(\mathbf{Z})$, the previous display is equal to $\mathbf{z}_k^\top \tilde{f}^{tt'}(\mathbf{Z}) - \sum_l \partial_{lk} \mathbf{e}_l^\top \tilde{f}^{tt'}(\mathbf{Z})$. We apply (69) again with respect to \mathbf{z}_k , so that

$$\begin{aligned} \mathbb{E} \text{Tr}[\mathbf{M}_j^\top \mathbf{M}_k] &= \sum_{il, tt'} \mathbf{e}_i^\top \partial_{lk} [f^{tt'}(\mathbf{Z})] \mathbf{e}_l^\top \tilde{f}^{tt'}(\mathbf{Z}) \\ &= \sum_{il, tt'} \left(\mathbf{e}_i^\top \partial_{lk} [\mathbb{E}_j \mathbf{V} e_t \mathbf{e}_{t'}^\top \mathbb{E}_j \mathbf{U}^\top] \mathbf{z}_j \right) \left(\mathbf{e}_l^\top \partial_{ij} [\mathbb{E}_k [\mathbf{U}] \mathbf{e}_{t'} \mathbf{e}_{t'}^\top \mathbb{E}_k [\mathbf{V}]^\top] \mathbf{z}_k \right). \end{aligned}$$

To remove the indices t, t' , we rewrite the above using $\sum_t \mathbf{e}_t \mathbf{e}_t^\top = \mathbf{I}_T$ and $\sum_{t'} \mathbf{e}_{t'} \mathbf{e}_{t'}^\top = \mathbf{I}_T$ so that it equals

$$\mathbb{E} \sum_{il} \text{Tr} \left\{ \partial_{lk} [\mathbb{E}_j \mathbf{U}^\top \mathbf{z}_j \mathbf{e}_i^\top \mathbb{E}_j \mathbf{V}] \partial_{ij} [\mathbb{E}_k [\mathbf{V}]^\top \mathbf{z}_k \mathbf{e}_l^\top \mathbb{E}_k [\mathbf{U}]] \right\}.$$

Summing over j, k , using $\text{Tr}[\mathbf{A}^\top \mathbf{B}] \leq \|\mathbf{A}\|_F \|\mathbf{B}\|_F$ and the Cauchy-Schwarz inequality, the above is bounded from above by

$$\left\{ \mathbb{E} \sum_{jk, il} \left\| \partial_{lk} [\mathbb{E}_j \mathbf{U}^\top \mathbf{z}_j \mathbf{e}_i^\top \mathbb{E}_j \mathbf{V}] \right\|_F^2 \right\}^{\frac{1}{2}} \left\{ \mathbb{E} \sum_{jk, il} \left\| \partial_{ij} [\mathbb{E}_k [\mathbf{V}]^\top \mathbf{z}_k \mathbf{e}_l^\top \mathbb{E}_k [\mathbf{U}]] \right\|_F^2 \right\}^{\frac{1}{2}}.$$

At this point the two factors are symmetric, with (\mathbf{V}, \mathbf{U}) in the left factor replaced with (\mathbf{U}, \mathbf{V}) on the right factor. We focus on the left factor; similar bound will apply to the right one by exchanging the

roles of \mathbf{V} and \mathbf{U} . If \mathbf{z}_j is independent of matrices $A^{(q)}$ $\mathbb{E}_j[\|\sum_{q=1}^n (\mathbf{e}_q^\top \mathbf{z}_j) A^{(q)}\|_F^2] = \sum_{q=1}^n \|A^{(q)}\|_F^2$ so that with $A^{(q)} = \partial_{lk}[\mathbb{E}_j \mathbf{U}^\top \mathbf{e}_q \mathbf{e}_i^\top \mathbb{E}_j \mathbf{V}]$, the first factor in the above display is equal to

$$\begin{aligned}
& \left\{ \mathbb{E} \sum_{jk, ilq} \left\| \partial_{lk} \left(\mathbb{E}_j \mathbf{U}^\top \mathbf{e}_q \mathbf{e}_i^\top \mathbb{E}_j \mathbf{V} \right) \right\|_F^2 \right\}^{\frac{1}{2}} \\
& \stackrel{(i)}{=} \left\{ \mathbb{E} \sum_{jk, ilq} \left\| \partial_{lk} \left(\mathbb{E}_j \mathbf{U}^\top \right) \mathbf{e}_q \mathbf{e}_i^\top \mathbb{E}_j [\mathbf{V}] + \mathbb{E}_j [\mathbf{U}]^\top \mathbf{e}_q \mathbf{e}_i^\top \partial_{lk} \left(\mathbb{E}_j [\mathbf{V}] \right) \right\|_F^2 \right\}^{\frac{1}{2}} \\
& \stackrel{(ii)}{\leq} \left\{ \mathbb{E} \sum_{jk, ilq} \left\| \partial_{lk} \left(\mathbb{E}_j [\mathbf{U}^\top] \right) \mathbf{e}_q \mathbf{e}_i^\top \mathbb{E}_j [\mathbf{V}] \right\|_F^2 \right\}^{\frac{1}{2}} + \left\{ \mathbb{E} \sum_{jk, ilq} \left\| \mathbb{E}_j \mathbf{U}^\top \mathbf{e}_q \mathbf{e}_i^\top \partial_{lk} \left(\mathbb{E}_j [\mathbf{V}] \right) \right\|_F^2 \right\}^{\frac{1}{2}} \\
& \stackrel{(iii)}{=} \left\{ \mathbb{E} \sum_{jk, ilq} \left\| \mathbb{E}_j [\partial_{lk} \mathbf{U}]^\top \mathbf{e}_q \right\|_2^2 \left\| \mathbf{e}_i^\top \mathbb{E}_j \mathbf{V} \right\|_2^2 \right\}^{\frac{1}{2}} + \left\{ \mathbb{E} \sum_{jk, ilq} \left\| \mathbb{E}_j \mathbf{U}^\top \mathbf{e}_q \right\|_2^2 \left\| \mathbf{e}_i^\top \mathbb{E}_j [\partial_{lk} \mathbf{V}] \right\|_2^2 \right\}^{\frac{1}{2}} \\
& \stackrel{(iv)}{=} \left\{ \mathbb{E} \sum_{jk, l} \left\| \mathbb{E}_j [\partial_{lk} \mathbf{U}]^\top \right\|_F^2 \left\| \mathbb{E}_j \mathbf{V} \right\|_F^2 \right\}^{\frac{1}{2}} + \left\{ \mathbb{E} \sum_{jk, l} \left\| \mathbb{E}_j \mathbf{U}^\top \right\|_F^2 \left\| \mathbb{E}_j [\partial_{lk} \mathbf{V}] \right\|_F^2 \right\}^{\frac{1}{2}},
\end{aligned}$$

where (i) is the chain rule, (ii) the triangle inequality, (iii) holds provided that the order of the derivation ∂_{lk} and the expectation sign \mathbb{E}_j can be switched and using $\|\mathbf{a}\mathbf{b}^\top\|_F^2 = \|\mathbf{a}\|_2^2 \|\mathbf{b}\|_2^2$ for vectors \mathbf{a}, \mathbf{b} , and (iv) holds using $\sum_i \|\mathbf{A}\mathbf{e}_i\|_2^2 = \|\mathbf{A}\|_F^2 = \sum_q \|\mathbf{A}\mathbf{e}_q\|_2^2$ for a matrix \mathbf{A} with n columns. Finally, by Jensen's inequality, the above display is bounded by

$$\left\{ \mathbb{E} \sum_{k, l} \left\| \partial_{lk} \mathbf{U} \right\|_F^2 \sum_j \left\| \mathbb{E}_j [\mathbf{V}] \right\|_F^2 \right\}^{\frac{1}{2}} + \left\{ \mathbb{E} \sum_{k, l} \left\| \partial_{lk} \mathbf{V} \right\|_F^2 \sum_j \left\| \mathbb{E}_j [\mathbf{U}] \right\|_F^2 \right\}^{\frac{1}{2}}.$$

Since $\|\mathbf{U}\|_F \vee \|\mathbf{V}\|_F \leq 1$ almost surely, the previous display is bounded by $\sqrt{p}(\|\mathbf{U}\|_\partial + \|\mathbf{V}\|_\partial)$. In summary,

$$\mathbb{E}[\|(66)\|_F^2]^{\frac{1}{2}} \leq (2p + [2\sqrt{p}(\|\mathbf{U}\|_\partial + \|\mathbf{V}\|_\partial)]^2)^{\frac{1}{2}} \leq \sqrt{2p} + 2\sqrt{p}(\|\mathbf{U}\|_\partial + \|\mathbf{V}\|_\partial).$$

Combining the bounds on the terms (65)-(66)-(67)-(68) with the triangle inequality completes the proof. ■

Proof of Lemma E.13. Since $\mathbf{M}^\dagger \preceq \mathbf{M}_1^\dagger = \mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + \tau n \mathbf{P}_{\hat{\mathcal{J}}})^\dagger$,

$$\begin{aligned}
\|\mathbf{N}\|_{\text{op}} &= \|(\mathbf{I}_T \otimes \mathbf{X}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \mathbf{X}^\top)\|_{\text{op}} \\
&\leq \|(\mathbf{I}_T \otimes \mathbf{X})(\mathbf{I}_T \otimes (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + \tau n \mathbf{P}_{\hat{\mathcal{J}}})^\dagger)(\mathbf{I}_T \otimes \mathbf{X}^\top)\|_{\text{op}} \\
&= \|\mathbf{X}_{\hat{\mathcal{J}}} (\mathbf{X}_{\hat{\mathcal{J}}}^\top \mathbf{X}_{\hat{\mathcal{J}}} + \tau n \mathbf{P}_{\hat{\mathcal{J}}})^\dagger \mathbf{X}_{\hat{\mathcal{J}}}^\top\|_{\text{op}} \\
&\leq 1,
\end{aligned}$$

where the first inequality uses $\|\mathbf{A}\mathbf{B}\mathbf{A}^\top\|_{\text{op}} \leq \|\mathbf{A}\mathbf{C}\mathbf{A}^\top\|_{\text{op}}$ for $0 \preceq \mathbf{B} \preceq \mathbf{C}$. ■

Proof of Lemma E.14. By Lemma E.4, $\frac{\partial F_{lt}}{\partial z_{ij}} = D_{ij}^{lt} + \Delta_{ij}^{lt}$, where

$$D_{ij}^{lt} = -(\mathbf{e}_j^\top \mathbf{H} \otimes \mathbf{e}_i^\top)(\mathbf{I}_{nT} - \mathbf{N})(\mathbf{e}_t \otimes \mathbf{e}_l), \quad (70)$$

$$\Delta_{ij}^{lt} = -(\mathbf{e}_t^\top \otimes \mathbf{e}_l^\top)(\mathbf{I}_T \otimes \mathbf{X})\mathbf{M}^\dagger(\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}})(\mathbf{F}^\top \otimes \mathbf{I}_p)(\mathbf{e}_i \otimes \mathbf{e}_j). \quad (71)$$

For the first equality, since $\mathbf{e}_i^\top \frac{\partial \mathbf{F}}{\partial z_{ij}} = \sum_t (D_{ij}^{it} + \Delta_{ij}^{it})\mathbf{e}_t^\top$, we have

$$\begin{aligned} & \sum_{i=1}^n \sum_{j=1}^p \mathbf{F}^\top \mathbf{Z} \mathbf{e}_j \mathbf{e}_i^\top \frac{\partial \mathbf{F}}{\partial z_{ij}} \\ &= \sum_{ij} \mathbf{F}^\top \mathbf{Z} \mathbf{e}_j \sum_{t=1}^T (D_{ij}^{it} + \Delta_{ij}^{it})\mathbf{e}_t^\top \\ &= \sum_{ij} \mathbf{F}^\top \mathbf{Z} \mathbf{e}_j \sum_t D_{ij}^{it} \mathbf{e}_t^\top + \underbrace{\sum_{ij} \mathbf{F}^\top \mathbf{Z} \mathbf{e}_j \sum_t \Delta_{ij}^{it} \mathbf{e}_t^\top}_{J_1}, \end{aligned}$$

where the first term can be simplified as below

$$\begin{aligned} & \sum_{j=1}^p \sum_{i=1}^n \mathbf{F}^\top \mathbf{Z} \mathbf{e}_j \sum_{t=1}^T D_{ij}^{it} \mathbf{e}_t^\top \\ &= - \sum_{j=1}^p \sum_{i=1}^n \mathbf{F}^\top \mathbf{Z} \mathbf{e}_j \sum_{t=1}^T (\mathbf{e}_j^\top \mathbf{H} \otimes \mathbf{e}_i^\top)(\mathbf{I}_{nT} - \mathbf{N})(\mathbf{e}_t \otimes \mathbf{e}_i) \mathbf{e}_t^\top \\ &= - \mathbf{F}^\top \mathbf{Z} \mathbf{H} \left[\sum_i (\mathbf{I}_T \otimes \mathbf{e}_i^\top)(\mathbf{I}_{nT} - \mathbf{N})(\mathbf{I}_T \otimes \mathbf{e}_i) \right] \\ &= - \mathbf{F}^\top \mathbf{Z} \mathbf{H} (n\mathbf{I}_T - \widehat{\mathbf{A}}). \end{aligned}$$

For the second equality, since $\frac{\partial \mathbf{F}}{\partial z_{ij}} = \sum_{lt} \mathbf{e}_l (D_{ij}^{lt} + \Delta_{ij}^{lt})\mathbf{e}_t^\top$,

$$\sum_{ij} \left(\frac{\partial \mathbf{F}}{\partial z_{ij}} \right)^\top \mathbf{Z} \mathbf{e}_j \mathbf{e}_i^\top \mathbf{F} = \underbrace{\sum_{ijtl} \mathbf{e}_t D_{ij}^{lt} \mathbf{e}_l^\top \mathbf{Z} \mathbf{e}_j \mathbf{e}_i^\top \mathbf{F}}_{J_2} + \sum_{ijtl} \mathbf{e}_t \Delta_{ij}^{lt} \mathbf{e}_l^\top \mathbf{Z} \mathbf{e}_j \mathbf{e}_i^\top \mathbf{F},$$

where the second term can be simplified as below,

$$\begin{aligned} & \sum_{ijtl} \mathbf{e}_t \Delta_{ij}^{lt} \mathbf{e}_l^\top \mathbf{Z} \mathbf{e}_j \mathbf{e}_i^\top \mathbf{F} \\ &= \sum_{ijtl} \mathbf{e}_t \Delta_{ij}^{lt} (\mathbf{e}_i \otimes \mathbf{e}_j)^\top (\mathbf{F} \otimes \mathbf{Z}^\top \mathbf{e}_l) \\ &= - \sum_{ijtl} \mathbf{e}_t (\mathbf{e}_t^\top \otimes \mathbf{e}_l^\top)(\mathbf{I}_T \otimes \mathbf{X})\mathbf{M}^\dagger(\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}})(\mathbf{F}^\top \otimes \mathbf{I}_p)(\mathbf{e}_i \otimes \mathbf{e}_j)(\mathbf{e}_i \otimes \mathbf{e}_j)^\top (\mathbf{F} \otimes \mathbf{Z}^\top \mathbf{e}_l) \\ &= - \sum_l (\mathbf{I}_T \otimes \mathbf{e}_l^\top)(\mathbf{I}_T \otimes \mathbf{X})\mathbf{M}^\dagger(\mathbf{I}_T \otimes \Sigma^{\frac{1}{2}})(\mathbf{F}^\top \otimes \mathbf{I}_p)(\mathbf{F} \otimes \mathbf{Z}^\top \mathbf{e}_l) \end{aligned}$$

$$\begin{aligned}
&= - \sum_l (\mathbf{I}_T \otimes \mathbf{e}_l^\top) (\mathbf{I}_T \otimes \mathbf{X}) \mathbf{M}^\dagger (\mathbf{F}^\top \mathbf{F} \otimes \mathbf{X}^\top \mathbf{e}_l) \\
&= - \sum_l (\mathbf{I}_T \otimes \mathbf{e}_l^\top) (\mathbf{I}_T \otimes \mathbf{X}) \mathbf{M}^\dagger (\mathbf{I}_T \otimes \mathbf{X}^\top \mathbf{e}_l) \mathbf{F}^\top \mathbf{F} \\
&= - \widehat{\mathbf{A}} \mathbf{F}^\top \mathbf{F},
\end{aligned}$$

where the last line uses the expression of $\widehat{\mathbf{A}}$ in (10).

It remains to bound the norm of \mathbf{J}_1 and \mathbf{J}_2 . To bound $\|\mathbf{J}_2\|_F$, recall the definition of \mathbf{J}_2 ,

$$\begin{aligned}
\mathbf{J}_2 &= \sum_{ijlt} \mathbf{e}_t D_{ij}^{lt} \mathbf{e}_l^\top \mathbf{Z} \mathbf{e}_j \mathbf{e}_i^\top \mathbf{F} \\
&= - \sum_{ijlt} \mathbf{e}_t (\mathbf{e}_j^\top \mathbf{H} \otimes \mathbf{e}_i^\top) (\mathbf{I}_{nT} - \mathbf{N}) (\mathbf{e}_t \otimes \mathbf{e}_l) \mathbf{e}_l^\top \mathbf{Z} \mathbf{e}_j \mathbf{e}_i^\top \mathbf{F} \\
&= - \sum_{ijlt} \mathbf{e}_t (\mathbf{e}_t^\top \otimes \mathbf{e}_l^\top) (\mathbf{I}_{nT} - \mathbf{N}) (\mathbf{H}^\top \mathbf{e}_j \otimes \mathbf{e}_i) \mathbf{e}_l^\top \mathbf{Z} \mathbf{e}_j \mathbf{e}_i^\top \mathbf{F} \\
&= - \sum_{ijl} (\mathbf{I}_T \otimes \mathbf{e}_l^\top) (\mathbf{I}_{nT} - \mathbf{N}) (\mathbf{H}^\top \mathbf{e}_j \otimes \mathbf{e}_i) \mathbf{e}_l^\top \mathbf{Z} \mathbf{e}_j \mathbf{e}_i^\top \mathbf{F} \\
&= - \sum_l (\mathbf{I}_T \otimes \mathbf{e}_l^\top) (\mathbf{I}_{nT} - \mathbf{N}) (\mathbf{H}^\top \mathbf{Z}^\top \mathbf{e}_l \otimes \mathbf{F}).
\end{aligned}$$

Since \mathbf{N} is non-negative definite with $\|\mathbf{N}\|_{\text{op}} \leq 1$, $\|\mathbf{I}_{nT} - \mathbf{N}\|_{\text{op}} \leq 1$,

$$\begin{aligned}
\|\mathbf{J}_2\|_F &\leq \sum_l \|(\mathbf{I}_T \otimes \mathbf{e}_l^\top) (\mathbf{I}_{nT} - \mathbf{N}) (\mathbf{H}^\top \mathbf{Z}^\top \mathbf{e}_l \otimes \mathbf{F})\|_F \\
&\leq \sum_l \|(\mathbf{I}_T \otimes \mathbf{e}_l^\top) (\mathbf{I}_{nT} - \mathbf{N})\|_{\text{op}} \|(\mathbf{H}^\top \mathbf{Z}^\top \mathbf{e}_l \otimes \mathbf{F})\|_F \\
&\leq \sum_l \|(\mathbf{H}^\top \mathbf{Z}^\top \mathbf{e}_l \otimes \mathbf{F})\|_F \\
&= \sum_l \|\mathbf{H}^\top \mathbf{Z}^\top \mathbf{e}_l\|_F \|\mathbf{F}\|_F \\
&\leq n^{\frac{1}{2}} \|\mathbf{Z} \mathbf{H}\|_F \|\mathbf{F}\|_F \\
&\leq n^{\frac{1}{2}} \|\mathbf{Z}\|_{\text{op}} \|\mathbf{H}\|_F \|\mathbf{F}\|_F.
\end{aligned}$$

To bound $\|\mathbf{J}_1\|_F$, recall the definition of \mathbf{J}_1 ,

$$\mathbf{J}_1 = \sum_{ij} \mathbf{F}^\top \mathbf{Z} \mathbf{e}_j \sum_t \Delta_{ij}^{it} \mathbf{e}_t^\top.$$

For each $t, t' \in [T]$,

$$\mathbf{e}_{t'}^\top \mathbf{J}_1 \mathbf{e}_t = \mathbf{e}_{t'}^\top \left[\sum_{ij} \mathbf{F}^\top \mathbf{Z} \mathbf{e}_j \sum_t \Delta_{ij}^{it} \mathbf{e}_t^\top \right] \mathbf{e}_t$$

$$\begin{aligned}
&= \sum_{j=1}^p \sum_{i=1}^n \mathbf{e}_{t'}^\top \mathbf{F}^\top \mathbf{Z} \mathbf{e}_j \Delta_{ij}^{it} \\
&= - \sum_i (\mathbf{e}_i^\top \mathbf{F} \otimes \mathbf{e}_{t'}^\top \mathbf{F}^\top) \mathbf{N}(\mathbf{e}_t \otimes \mathbf{e}_i) \\
&= -\mathbf{e}_{t'}^\top \sum_i (\mathbf{e}_i^\top \mathbf{F} \otimes \mathbf{F}^\top) \mathbf{N}(\mathbf{I}_T \otimes \mathbf{e}_i) \mathbf{e}_t.
\end{aligned}$$

Thus, $\mathbf{J}_1 = \sum_i (\mathbf{e}_i^\top \mathbf{F} \otimes \mathbf{F}^\top) \mathbf{N}(\mathbf{I}_T \otimes \mathbf{e}_i)$, and hence

$$\begin{aligned}
\|\mathbf{J}_1\|_F &= \left\| \sum_i (\mathbf{e}_i^\top \mathbf{F} \otimes \mathbf{F}^\top) \mathbf{N}(\mathbf{I}_T \otimes \mathbf{e}_i) \right\|_F \\
&\leq \sum_i \|(\mathbf{e}_i^\top \mathbf{F} \otimes \mathbf{F}^\top) \mathbf{N}(\mathbf{I}_T \otimes \mathbf{e}_i)\|_F \\
&\leq \sum_i \|(\mathbf{e}_i^\top \mathbf{F} \otimes \mathbf{F}^\top)\|_F \|\mathbf{N}(\mathbf{I}_T \otimes \mathbf{e}_i)\|_{\text{op}} \\
&\leq \sum_i \|\mathbf{e}_i^\top \mathbf{F}\| \|\mathbf{F}\|_F \\
&\leq n^{\frac{1}{2}} \|\mathbf{F}\|_F^2,
\end{aligned}$$

where the first inequality is by sub-additivity of Frobenius norm, the second inequality uses $\|\mathbf{A}_1 \mathbf{A}_2\|_F \leq \|\mathbf{A}_1\|_F \|\mathbf{A}_2\|_{\text{op}}$ for any matrices $\mathbf{A}_1, \mathbf{A}_2$ with appropriate dimensions, the third inequality is by $\|\mathbf{N}\|_{\text{op}} \leq 1$ from Lemma E.13, the last inequality is by Cauchy-Schwarz inequality. ■