

Supplementary materials for the paper “Nonparametric Eigenvalue-Regularized Precision or Covariance Matrix Estimator”- Simulation results and the proof of Theorem 1, Theorem 3, Lemma 1, Theorem 5 and Theorem 6 in the paper.

1 Simulation Experiments

We demonstrate and compare the performance of our estimator to other state-of-the-art estimators under various settings. Hereafter, we abbreviate our method as NERCOME for estimating Σ_p or Ω_p , which comes from the name Nonparametric Eigenvalue-Regularized COvariance Matrix Estimator. The method proposed in Abadir et al. (2014) is abbreviated as CRC (Condition number Regularized Covariance estimator), while the nonlinear shrinkage method in Ledoit and Wolf (2012) is abbreviated as NONLIN. We call the grand average estimator (15) in Abadir et al. (2014) the CRC grand average. Finally, the method in Fan et al. (2013) is abbreviated as POET. We also include in some cases the graphical LASSO, abbreviated as GLASSO, as in Friedman et al. (2008), and the adaptive SCAD thresholding, abbreviated as SCAD, which is a special case of POET without any factors. We create five different profiles:

- (I) Independent data $\mathbf{y}_t \sim N(\mathbf{0}, \Sigma_p)$, where $\Sigma_p = \mathbf{QDQ}^T$. The orthogonal matrix \mathbf{Q} is randomly generated each time, and \mathbf{D} is diagonal with 40% of values being 3, and 60% being 7.
- (II) Same as (I), except that $\Sigma_p = \mathbf{I}_p$.
- (III) The data is from a factor model $\mathbf{y}_t = \mathbf{A}\mathbf{x}_t + \epsilon_t$, where \mathbf{A} has size $p \times 3$, with elements in \mathbf{A} generated independently from $N(0, 2^2)$. The \mathbf{x}_t 's are independently generated from $N(\mathbf{0}, 2\mathbf{I}_r)$, while the ϵ_t 's are independently generated from the standardized t_3 distribution.
- (IV) (Sparse covariance matrix). The covariance matrix Σ_p is sparse and is randomly generated each time, with 20% non-zeros, and independent data $\mathbf{y}_t \sim N(\mathbf{0}, \Sigma_p)$.
- (V) Same as (IV), except that $\mathbf{y}_t = \Sigma_p \mathbf{z}_t$ where the \mathbf{z}_t 's are independent containing independent standardized t_5 random variables.

Profiles (I), (II), (IV) and (V) are for non-factor model, with profile (II) being excluded in Theorem 5, since the loss of an ideal estimator is exactly 0 when $\Sigma_p = \sigma^2 \mathbf{I}_p$. Profile (III) is for factor model. Profiles (III) and (V) test the robustness of our method to fat-tailed distributions, since assumption (F1) rules out an error distribution of t_3 , while assumption (A1) rules out t_5 being the distribution of the independent random variables.

Recall that NONLIN is proved to be asymptotically optimal with respect to the Frobenius and the inverse Stein's loss functions when the data is not from a factor model (see the results in Ledoit and Wolf (2013a) for more details).

Hence we expect that NERCOME should be close to NONLIN in performance for the non-factor model profiles, but better than NONLIN for profile (III), since NERCOME is proved to be asymptotically optimal in Theorem 3 for data from a factor model.

We simulate 500 times from the seven profiles under all different combinations of $n = 200, 400, 800$ and $p = 50, 100, 200, 500$, and calculate the mean efficiency loss defined in (4.2) with respect to the Frobenius loss and the inverse Stein's loss given in (4.3) and (2.16) respectively. For profiles (IV) and (V), we include 5 more loss functions for comparisons. See Table 1 and the paragraphs therein for details. For both NERCOME and CRC, we use $M = 50$ permutations for each split of the data, which gives a good trade-off between accuracy and computational complexity. We use the 7 different split locations as in (4.8) for NERCOME and CRC. We also demonstrate the performance of NERCOME using the automatic choice of m by minimizing $g(m)$ defined in (4.7). The CRC grand average estimator defined in (18) of Abadir et al. (2014) will also be included for comparison, which takes average over the 4 CRC estimators calculated on the split locations $0.2n, 0.4n, 0.6n$ and $0.8n$ respectively.

In the titles of the subsequent graphs, the inverse Stein's loss is abbreviated as the Stein's loss to save space.

Figure 1 shows the mean efficiency loss over 500 simulations for various methods under profile (I). Clearly, the graphical LASSO and the SCAD adaptive thresholding are not performing as good as other methods in terms of the Frobenius loss function, since the covariance and the precision matrices under profile (I) are only approximately sparse. The graphical LASSO has the worst performance, which is to be expected since it concerns with the sparse precision matrix rather than the covariance matrix estimation. NONLIN is the best in all scenarios, followed by NERCOME with m chosen automatically by minimizing (4.7), and CRC grand average. It is fair to say that NONLIN, NERCOME and CRC grand average all perform well in absolute terms. Comparing the performance at different split locations for NERCOME and CRC, we can see that CRC attains a minimum with a smaller split m , and performs better when m is not large. The opposite is true for NERCOME, which usually is at the smallest efficiency loss when m is large. This is consistent with the condition $m/n \rightarrow 1$ for asymptotic efficiency in Theorem 5. Clearly, the choice of m that minimizes $g(m)$ in (4.7) is good, since the resulting mean efficiency loss is close to the minimum for NERCOME. We can also see that for CRC, averaging over estimators with split locations from $0.2n$ to $0.8n$, which is the CRC grand average estimator suggested in Abadir et al. (2014), may result in suboptimal performance, since the CRC is usually suboptimal when m is small.

We have omitted the mean efficiency loss with respect to the inverse Stein's loss for profile (I), since the graphs are very similar to the Frobenius loss's counterparts in Figure 1.

Figure 2 shows the actual Frobenius loss and inverse Stein's loss for various methods for profile (II). The efficiency loss in (4.2) is always 1 for any imperfect estimators, since it is easy to see that $\hat{\Sigma}_{\text{Ideal}} = \Sigma_p = \sigma^2 \mathbf{I}_p$. Even for $\Sigma_p = \sigma^2 \mathbf{I}_p$,

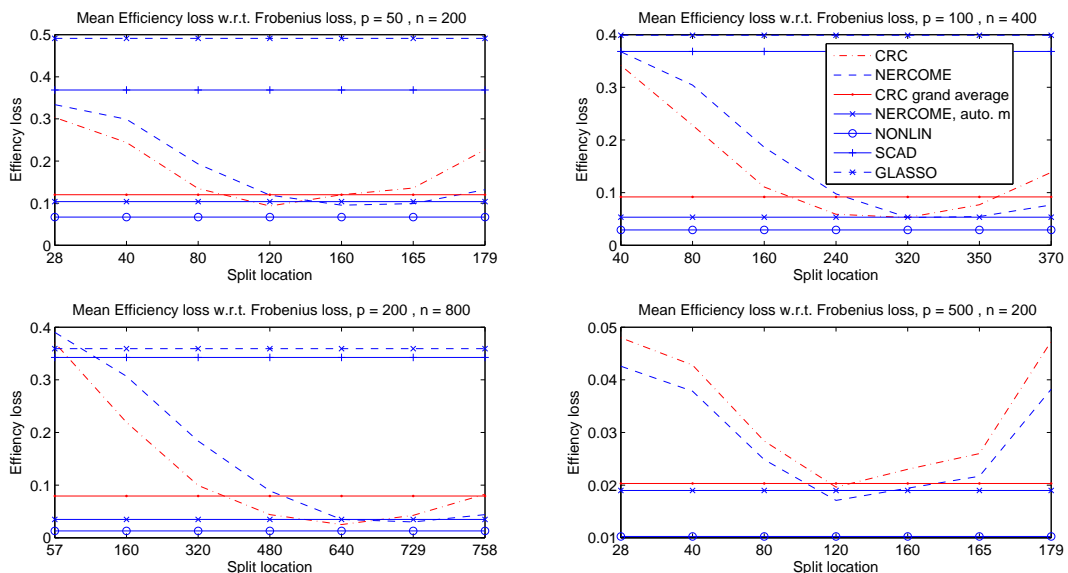


Figure 1: Mean efficiency loss with respect to the Frobenius loss for profile (I). SCAD and GLASSO are included only when $p < 500$ to save computational time. NERCOME with m automatically chosen is to minimize $g(m)$ in (4.7) over the split locations in (4.8). CRC shares the same split locations. Only NERCOME and CRC are varying over different split locations.

SCAD thresholding and graphical LASSO are not as good as other methods in terms of both loss functions. This can be explained by the fact that even if all the off-diagonal elements of the covariance and the precision matrix estimated from the data are correctly killed off, the diagonal elements are left untouched. Hence, the errors in the diagonal elements will still stack up in the Frobenius or the inverse Stein's loss. On the other hand, CRC, NERCOME and NONLIN are in principal setting the eigenvalues to be as close to $\mathbf{p}_i^T \boldsymbol{\Sigma}_p \mathbf{p}_i$ as possible. Under $\boldsymbol{\Sigma}_p = \sigma^2 \mathbf{I}_p$, $\mathbf{p}_i^T \boldsymbol{\Sigma}_p \mathbf{p}_i = \sigma^2$, which is the correct magnitude. The off-diagonal entries will then be close to 0 since $\mathbf{P} \text{diag}(\sigma^2, \dots, \sigma^2) \mathbf{P}^T = \mathbf{P}(\sigma^2 \mathbf{I}_p) \mathbf{P}^T = \sigma^2 \mathbf{I}_p$, no matter what orthogonal \mathbf{P} we use. Hence, we can see that the advantage of CRC, NERCOME and NONLIN comes from the fact that the class of estimators $\boldsymbol{\Sigma}(\mathbf{D}) = \mathbf{P} \mathbf{D} \mathbf{P}^T$ is particularly good when the true covariance matrix is diagonal, with the same entries on the main diagonal throughout. The actual performance of CRC, NERCOME and NONLIN are very close in this case.

It is clear that for both NERCOME and CRC, the best split location is when m is the smallest. This can be seen for NERCOME by noting that, when $\boldsymbol{\Sigma}_p = \sigma^2 \mathbf{I}_p$, $\|\hat{\boldsymbol{\Sigma}}_m - \boldsymbol{\Sigma}_p\|_F^2 = \sum_{i=1}^p (\mathbf{p}_{1i}^T \tilde{\boldsymbol{\Sigma}}_2 \mathbf{p}_{1i} - \sigma^2)^2$. Hence, in order to minimize

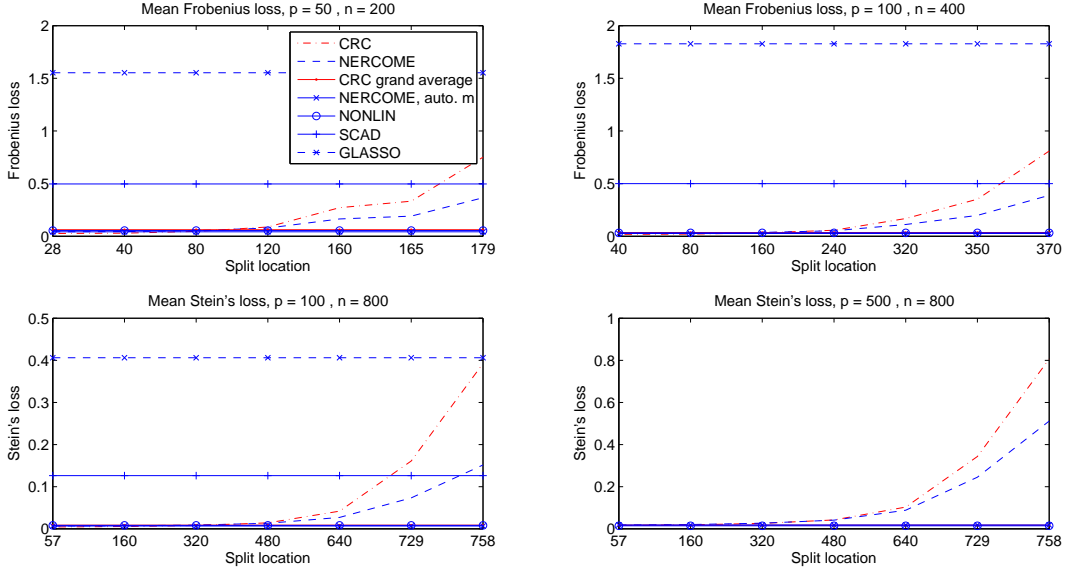


Figure 2: Mean Frobenius loss (upper row) and inverse Stein's loss (lower row) for profile (II). SCAD and GLASSO are absent for $p < 500$. Refer to the descriptions of Figure 1 for further details.

the loss, we want $\tilde{\Sigma}_2$ to be as close to $\sigma^2 \mathbf{I}_p$ as possible. It means that $n_2 = n - m$ should be made as large as possible, which implies that m should be as small as possible. The automatic selection for NERCOME clearly performs well even in this case.

For profile (III), all the methods perform badly in terms of Frobenius loss with high efficiency losses, hence the corresponding figure is not shown. This can be explained by the fact that since the factor model has 3 spiked eigenvalues of order p , any error in estimating those eigenvalues will be hugely reflected in the efficiency loss.

On the other hand, Figure 3 shows a huge improvement for CRC and NERCOME when the inverse Stein's loss is concerned, while POET performs the worst and NONLIN stays at around 50% efficiency loss. In one simulation run for $n = p = 200$, NONLIN returns a singular covariance matrix, and hence the inverse Stein's loss is undefined. This case has to be removed for calculating the mean of efficiency loss. Note that NONLIN always results in a positive semi-definite covariance matrix. However, when $n = p = 200$, it corresponds to $c = 1$ with exactly 1 zero eigenvalue for the sample covariance matrix. This case is not covered in Theorem 4 of Ledoit and P  ch   (2011), and the fact that the data is from a factor model violates the important assumption in Ledoit and P  ch  

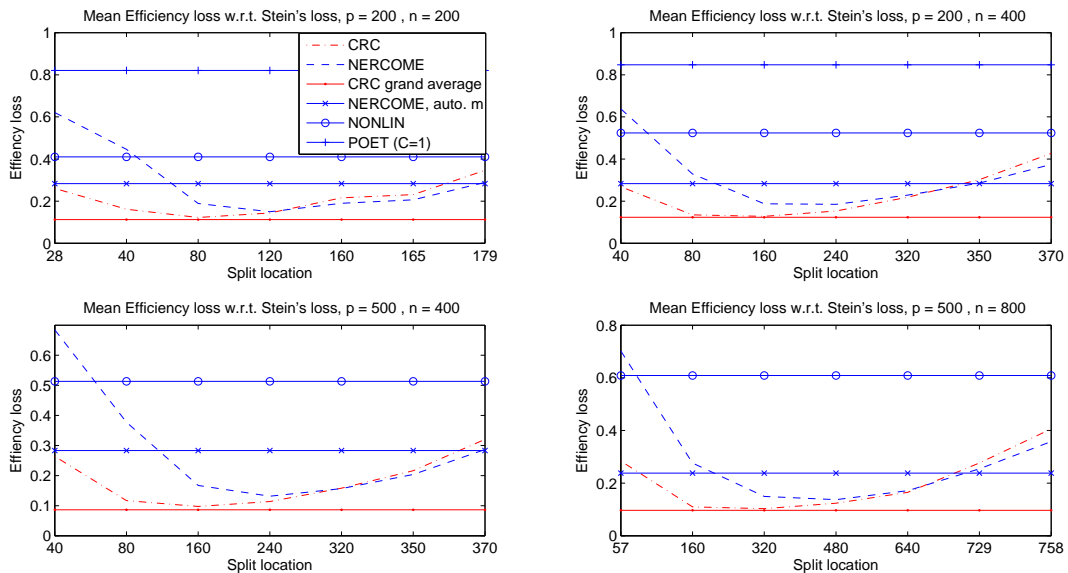


Figure 3: Mean efficiency loss with respect to the inverse Stein's loss for profile (III). The mean efficiency for NONLIN when $n = p = 200$ (upper left panel) is calculated after removing one simulation result with singular covariance matrix.

(2011) and Ledoit and Wolf (2013a,b), that we are able to write $\mathbf{y}_t = \boldsymbol{\Sigma}_p^{1/2} \mathbf{z}_t$, with \mathbf{z}_t having independent and identically distributed entries. Hence, although rare in practice, it is not surprising to see a singular covariance matrix estimator for NONLIN. On the other hand, we have proved in Corollary 4 that our estimator is almost surely positive definite even for data from a factor model, although assumption (F1) is not satisfied here.

CRC grand average performs the best now, since it appears that all the split locations between $0.2n$ and $0.8n$ are indeed giving good performance to CRC and NERCOME. Part of the reason that POET is not as good as CRC, NERCOME and NONLIN is related to the fact that $\boldsymbol{\Sigma}_\epsilon = \mathbf{I}_p$, which gives advantages to the class of estimators $\boldsymbol{\Sigma}(\mathbf{D}) = \mathbf{PDP}^T$, with arguments similar to the explanations of why SCAD thresholding and graphical LASSO perform worse under profile (II).

For profiles (IV) and (V), we investigate the performance of various estimators in estimating sparse covariance matrix. NONLIN, NERCOME and CRC are all rotation-equivariant estimators, and regularize through stabilizing the eigenvalues. Hence they are not recovering the sparse structure, unlike the thresholding estimator SCAD. On top of the two original loss functions, we

introduce 5 more to gauge the performance:

$$\begin{aligned}
L_3(\boldsymbol{\Sigma}, \widehat{\boldsymbol{\Sigma}}) &= \|\boldsymbol{\Sigma}^{-1} - \widehat{\boldsymbol{\Sigma}}^{-1}\|_F, & L_4(\boldsymbol{\Sigma}, \widehat{\boldsymbol{\Sigma}}) &= \text{tr}(\boldsymbol{\Sigma}^{-1}\widehat{\boldsymbol{\Sigma}}) - \log \det(\boldsymbol{\Sigma}^{-1}\widehat{\boldsymbol{\Sigma}}) - p, \\
L_5(\boldsymbol{\Sigma}, \widehat{\boldsymbol{\Sigma}}) &= \|\boldsymbol{\Sigma} - \widehat{\boldsymbol{\Sigma}}\|, & L_6(\boldsymbol{\Sigma}, \widehat{\boldsymbol{\Sigma}}) &= \sum_{i=1}^p |\lambda_i(\boldsymbol{\Sigma} - \widehat{\boldsymbol{\Sigma}})|, \\
L_7(\boldsymbol{\Sigma}, \widehat{\boldsymbol{\Sigma}}) &= \text{tr}(\boldsymbol{\Sigma} + \widehat{\boldsymbol{\Sigma}} - 2\boldsymbol{\Sigma}^{1/2}\widehat{\boldsymbol{\Sigma}}^{1/2}).
\end{aligned}$$

L_3 is the Frobenius norm on the difference of the inverse covariance matrices; L_4 is the Stein's loss; L_5 is the operator/spectral norm of the difference; L_6 is the nuclear norm of the difference. Finally, L_7 is called the Fréchet loss. We use L_1 and L_2 for the Frobenius and the inverse Stein's loss, respectively.

Table 1 shows the results for profile (IV) when $p = 200$. Other combinations of (n, p) are omitted to save space. Sample covariance always leads to huge efficiency loss and hence is not shown. It is clear that when $\boldsymbol{\Sigma}_p$ is sparse, SCAD thresholding can take advantage and usually has negative efficiency loss when the dimension is not too close to the sample size. However, if the dimension is close to the sample size, SCAD has difficulty in finding a good thresholding parameter, resulting in large efficiency losses across different loss functions. NONLIN usually outperforms CRC grand average and NERCOME apart from SCAD, and efficiency loss is usually getting smaller for NONLIN when sample size gets larger while dimension stays fixed. NERCOME follows NONLIN closely in all (p, n) combinations, and outperforms in several occasions, especially when dimension is close to the sample size. In fact NERCOME seems to do particularly well in the Frobenius loss for the inverse and the Stein's loss. CRC grand average, on the other hand, does not usually have decreasing efficiency loss as sample size increases. On the contrary, many actually increase as sample size increases while dimension stays fixed, showing that averaging on different splits is not always a good idea, especially when n is large compared to p , which is a setting that Abadir et al. (2014) used. Figure 1 also shows this phenomenon. In all cases, we can see that among the three eigenvalues-stabilizing methods, NERCOME either outperforms, or is close to the best performer.

The same table shows the results for profile (V) when $p \geq 100$. With fat-tailed distribution, it is more difficult for SCAD to find a good thresholding parameter (except when n is 8 or more times larger than p). CRC grand average still shows an increasing efficiency loss in general when n is getting larger. NONLIN is now less dominant among the three eigenvalues-stabilizing methods, with NERCOME outperforming both of the other methods in more (p, n) combinations. NONLIN has produced two singular covariance matrices in the $p = n = 200$ simulations, showing that when assumption (A1) is violated, NONLIN can in fact produce singular covariance matrix. NERCOME and CRC grand average still produce positive definite covariance matrices in those two cases. Again, among the three methods, NERCOME either outperforms or is close to the best performer in all (p, n) settings.

Figure 4 shows the boxplots of the efficiency loss of the averaged ideal estimator $\widehat{\boldsymbol{\Sigma}}_{\text{Ideal}, m, M} = M^{-1} \sum_{i=1}^M \widehat{\boldsymbol{\Sigma}}_{\text{Ideal}, m}^{(i)}$ for various profiles. We show this

		Profile (IV)						
		L_1	L_2	L_3	L_4	L_5	L_6	L_7
$p = 200,$ $n = 200$	NERCOME	3.8 _(1.0)	3.5 _(1.0)	0.7 _(1.3)	1.0 _(1.6)	5.2 _(1.8)	1.7 _(0.4)	2.5 _(0.6)
	CRC	3.0 _(0.5)	2.7 _(0.5)	1.8 _(1.1)	1.7 _(1.4)	2.8 _(1.3)	1.3 _(0.3)	2.5 _(0.5)
	NONLIN	2.6 _(1.8)	3.5 _(6.7)	2.4 _(10.6)	1.2 _(3.2)	3.6 _(6.7)	1.2 _(0.7)	2.3 _(2.6)
	SCAD	18.2 _(1.0)	16.7 _(0.9)	12.2 _(1.0)	16.4 _(1.3)	10.4 _(2.3)	9.5 _(0.5)	17.6 _(1.0)
$p = 200,$ $n = 400$	NERCOME	2.8 _(0.5)	2.5 _(0.4)	1.2 _(1.1)	0.9 _(1.2)	3.2 _(1.2)	1.2 _(0.2)	2.2 _(0.4)
	CRC	3.6 _(0.5)	3.3 _(0.5)	4.7 _(0.9)	4.5 _(1.1)	2.1 _(1.2)	1.5 _(0.3)	3.7 _(0.6)
	NONLIN	1.7 _(0.4)	1.7 _(0.5)	0.6 _(0.9)	0.8 _(1.0)	1.6 _(1.3)	0.8 _(0.2)	1.5 _(0.4)
	SCAD	18.8 _(0.7)	17.0 _(0.7)	15.7 _(0.9)	19.4 _(1.2)	8.6 _(2.3)	9.6 _(0.4)	18.6 _(0.8)
$p = 200,$ $n = 800$	NERCOME	2.6 _(0.3)	2.5 _(0.3)	4.0 _(1.0)	3.1 _(0.9)	1.9 _(1.0)	1.1 _(0.2)	2.6 _(0.4)
	CRC	4.6 _(0.5)	4.9 _(0.5)	9.5 _(1.0)	8.1 _(1.0)	1.7 _(1.2)	2.1 _(0.2)	5.4 _(0.6)
	NONLIN	1.5 _(0.3)	1.5 _(0.3)	1.0 _(0.8)	1.0 _(0.8)	1.2 _(0.9)	0.7 _(0.2)	1.4 _(0.3)
	SCAD	6.6 _(1.4)	5.6 _(1.4)	5.5 _(1.6)	5.7 _(1.6)	6.0 _(2.8)	3.6 _(0.7)	6.1 _(1.4)
		Profile (V)						
$p = 100,$ $n = 200$	NERCOME	7.2 _(2.4)	6.9 _(2.5)	0.5 _(2.6)	0.3 _(3.2)	8.0 _(3.4)	3.2 _(1.1)	5.5 _(1.9)
	CRC	5.7 _(2.2)	5.2 _(1.9)	5.1 _(3.2)	5.1 _(4.1)	4.7 _(3.1)	2.4 _(1.1)	5.4 _(2.4)
	NONLIN	8.1 _(10.5)	6.4 _(3.4)	1.0 _(3.1)	2.9 _(7.1)	9.5 _(15.1)	3.7 _(3.5)	6.4 _(7.8)
	SCAD	39.0 _(7.6)	34.3 _(2.8)	25.4 _(2.4)	34.2 _(5.5)	33.2 _(15.6)	20.7 _(3.1)	36.7 _(5.7)
$p = 100,$ $n = 400$	NERCOME	5.0 _(1.4)	4.7 _(1.3)	3.4 _(2.7)	2.8 _(2.6)	5.2 _(2.7)	2.2 _(0.6)	4.3 _(1.1)
	CRC	6.2 _(1.1)	6.4 _(1.1)	10.6 _(2.6)	9.2 _(2.7)	3.8 _(2.6)	2.7 _(0.5)	6.8 _(1.2)
	NONLIN	5.1 _(4.0)	4.4 _(1.5)	1.0 _(2.0)	2.1 _(2.9)	5.4 _(10.7)	2.4 _(1.1)	4.2 _(2.6)
	SCAD	11.7 _(6.8)	10.6 _(4.0)	3.3 _(4.0)	7.8 _(5.6)	14.9 _(16.3)	4.9 _(2.5)	10.5 _(5.4)
$p = 100,$ $n = 800$	NERCOME	4.3 _(0.8)	4.3 _(0.7)	7.0 _(2.3)	4.9 _(2.0)	3.6 _(2.6)	1.9 _(0.4)	4.3 _(0.7)
	CRC	7.1 _(1.0)	8.4 _(1.0)	17.4 _(2.6)	12.5 _(2.1)	4.1 _(2.7)	3.3 _(0.5)	8.4 _(1.0)
	NONLIN	4.4 _(4.2)	3.9 _(1.4)	1.5 _(1.7)	2.4 _(2.7)	4.0 _(9.4)	2.0 _(1.0)	3.8 _(2.7)
	SCAD	9.9 _(8.0)	8.7 _(4.7)	16.3 _(5.1)	14.5 _(6.8)	1.0 _(17.4)	6.0 _(2.5)	10.8 _(6.4)
$p = 200,$ $n = 200$	NERCOME	4.3 _(1.4)	4.0 _(1.5)	0.9 _(1.8)	1.3 _(2.2)	5.7 _(2.4)	1.9 _(0.6)	2.8 _(0.9)
	CRC	3.2 _(0.8)	2.8 _(0.7)	1.9 _(1.8)	1.9 _(2.3)	3.3 _(1.8)	1.4 _(0.3)	2.7 _(0.8)
	*NONLIN	5.0 _(6.8)	4.0 _(5.9)	1.0 _(8.7)	1.3 _(4.5)	9.3 _(16.6)	2.1 _(1.6)	3.7 _(4.7)
	SCAD	34.3 _(6.3)	29.1 _(2.2)	19.0 _(1.8)	27.5 _(4.1)	41.6 _(17.5)	17.4 _(2.0)	31.2 _(4.2)
$p = 200,$ $n = 400$	NERCOME	2.9 _(0.6)	2.6 _(0.6)	1.5 _(1.8)	1.3 _(2.0)	3.0 _(1.8)	1.3 _(0.3)	2.4 _(0.7)
	CRC	3.5 _(0.6)	3.3 _(0.5)	4.7 _(1.4)	4.6 _(1.8)	1.9 _(1.8)	1.5 _(0.3)	3.7 _(0.6)
	NONLIN	3.5 _(4.9)	2.5 _(1.1)	0.0 _(1.4)	0.9 _(2.6)	6.3 _(14.6)	1.5 _(0.9)	2.5 _(2.7)
	SCAD	29.7 _(4.4)	25.6 _(1.5)	20.3 _(1.3)	27.2 _(2.7)	29.6 _(16.4)	15.2 _(1.2)	28.1 _(2.8)
$p = 200,$ $n = 800$	NERCOME	2.6 _(0.4)	2.5 _(0.3)	4.0 _(1.3)	3.1 _(1.4)	1.8 _(1.5)	1.1 _(0.2)	2.6 _(0.4)
	CRC	4.6 _(0.5)	4.9 _(0.5)	9.4 _(1.3)	8.1 _(1.4)	1.6 _(1.6)	2.0 _(0.2)	5.4 _(0.6)
	NONLIN	2.5 _(4.4)	1.9 _(0.8)	0.1 _(1.1)	0.8 _(2.2)	3.5 _(10.6)	1.1 _(0.7)	1.9 _(2.4)
	SCAD	7.2 _(5.2)	5.8 _(1.9)	5.7 _(2.0)	8.0 _(3.3)	11.8 _(16.9)	2.7 _(1.2)	7.0 _(3.4)

Table 1: Mean efficiency loss (%) with respect to different loss functions for NERCOME, CRC grand average NONLIN and SCAD for profile (IV) and (V). The standard deviation is in bracket (%). L_1 =Frobenius loss, L_2 =inverse Stein's loss, L_3 =Frobenius loss of inverse, L_4 =Stein's loss, L_5 =Spectral norm of difference, L_6 =nuclear norm of difference, L_7 =Fréchet loss. Shaded cells mean the number is negative. Bolded cells indicate the minimum among all methods. For profile (V), when $p = n = 200$, NONLIN has two cases of singular covariance matrix, which has to be removed for mean and standard deviation calculations.

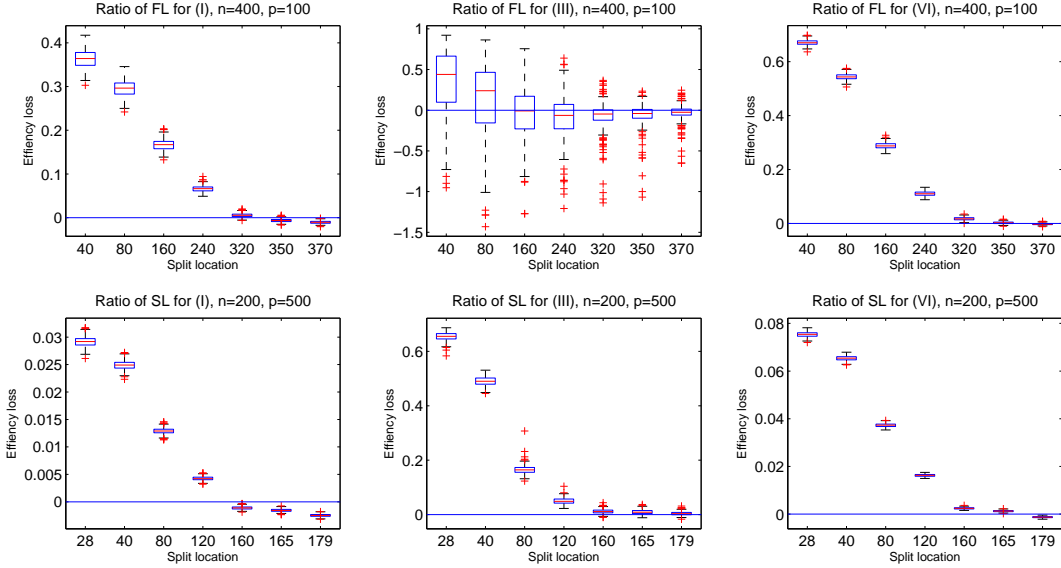


Figure 4: Boxplots of efficiency loss of $M^{-1} \sum_{i=1}^M \widehat{\Sigma}_{Ideal,m}^{(i)}$ at split locations defined in (4.8). Upper row: With respect to the Frobenius loss, with $n = 400, p = 100$. Lower row: With respect to the inverse Stein's loss, with $n = 200, p = 500$. Left panel : For profile (I). Right panel: For profile (III).

because it is important for the ratio of loss $L(\Sigma_p, \widehat{\Sigma}_{Ideal,m,M})/L(\Sigma_p, \widehat{\Sigma}_{Ideal})$ to be going to 1 almost surely (i.e., $EL(\Sigma_p, \widehat{\Sigma}_{Ideal,m,M}) \xrightarrow{a.s.} 0$, which is what Figure 4 is showing) to ensure our estimator $\widehat{\Sigma}_{m,M}$ to be almost surely asymptotically efficient. From the figure, it is clear that for the profiles (I) and (III), the split location m should be close to n for $EL(\Sigma_p, \widehat{\Sigma}_{Ideal,m,M})$ to be close to 0. The efficiency loss has a wide variability for profile (III) with respect to the Frobenius loss (the upper right panel of Figure 4). It certainly achieves 0 on average for several split locations not close to n , so that the best m is not always close to n for factor model with fat-tailed error distribution.

For profile (I), with negative efficiency loss in almost all 500 simulation runs for the averaged ideal estimator $\widehat{\Sigma}_{Ideal,m,M}$ when split location is close to n (see the left panels of Figure 4), it is clear that using the eigenmatrix \mathbf{P}_1 to form our estimator and averaging over different permutations of the data can be better than just using \mathbf{P} from the full data set. Our estimator $\widehat{\Sigma}_{m,M}$ in (4.6) follows the same spirit exactly.

2 Proof of theorems

We state an important Lemma first, which is in fact Lemma 2.7 of Bai and Silverstein (1998). Then, we present two more lemmas before presenting the proof of Theorem 1.

Lemma S.1 *Let $\mathbf{y} = (y_1, \dots, y_p)^\top$ be a complex random vector with independent and identically distributed (i.i.d.) entries satisfying $E(y_1) = 0$, $E|y_1|^2 = 1$. Let \mathbf{A} be a given $p \times p$ matrix with possibly complex entries. Then for any $q \geq 2$,*

$$E|\mathbf{y}^\top \mathbf{A} \mathbf{y} - \text{tr}(\mathbf{A})|^q \leq K_q \left(E^{q/2} |y_1|^{4\text{tr}^{q/2}}(\mathbf{A} \mathbf{A}^*) + E|y_1|^{2q} \text{tr}(\mathbf{A} \mathbf{A}^*)^{q/2} \right),$$

where K_q is a constant depending only on q , and $\mathbf{A}^* = \bar{\mathbf{A}}^\top$ is the Hermitian of \mathbf{A} .

Lemma S.2 *For $z \in \mathbb{C}^+$, we have $\|R_1(z)\|^2 \leq \frac{1}{\text{Im}^2(z)}$.*

Proof of Lemma S.2. Denote \bar{z} the conjugate of z , and let $\tilde{\Sigma}_1 = \mathbf{P}_1 \mathbf{D}_1 \mathbf{P}_1^\top$. Then

$$\begin{aligned} \|R_1(z)\|^2 &= \lambda_{\max}(R_1(z) \overline{R_1(z)}) \\ &= \lambda_{\max}(\mathbf{P}_1 (\mathbf{D}_1 - z \mathbf{I}_p)^{-1} \mathbf{P}_1^\top \mathbf{P}_1 (\mathbf{D}_1 - \bar{z} \mathbf{I}_p)^{-1} \mathbf{P}_1^\top) \\ &= \lambda_{\max}(\{[\mathbf{D}_1 - \text{Re}(z)]^2 + \text{Im}^2(z)\}^{-1}) \\ &\leq \frac{1}{\text{Im}^2(z)}. \quad \square \end{aligned}$$

Lemma S.3 *Let \mathbf{R} be Hermitian, that is, it is symmetric with complex entries. Then for a real square matrix \mathbf{B} , we have $|\text{tr}(\mathbf{B} \mathbf{R})| \leq |\lambda_{\max}(\mathbf{B})| |\text{tr}(\mathbf{R})|$. Also, $|\text{tr}(\mathbf{A}^\top \mathbf{R} \mathbf{A})| \leq 2^{1/2} \|\mathbf{R}\| |\text{tr}(\mathbf{A}^\top \mathbf{A})|$ for a real matrix \mathbf{A} .*

Proof of Lemma S.3. Since \mathbf{R} is Hermitian, it can be decomposed as $\mathbf{R} = \mathbf{U} \mathbf{D} \mathbf{U}^\top$, where \mathbf{U} is orthogonal and \mathbf{D} is diagonal with complex entries. Writing $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_p)$ and $\mathbf{D} = \text{diag}(d_1, \dots, d_p)$, then

$$\begin{aligned} |\text{tr}(\mathbf{B} \mathbf{R})| &= |\text{tr}(\mathbf{U}^\top \mathbf{B} \mathbf{U} \mathbf{D})| = \left| \sum_{i=1}^p d_i \mathbf{u}_i^\top \mathbf{B} \mathbf{u}_i \right| \\ &\leq |\lambda_{\max}(\mathbf{B})| \cdot \left| \sum_{i=1}^p d_i \right| = |\lambda_{\max}(\mathbf{B})| |\text{tr}(\mathbf{R})|. \end{aligned}$$

Finally, writing $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_r)$,

$$\begin{aligned} |\text{tr}(\mathbf{A}^\top \mathbf{R} \mathbf{A})| &= \left| \sum_{i=1}^r \mathbf{a}_i^\top \mathbf{R} \mathbf{a}_i \right| = \left(\left[\sum_{i=1}^r \mathbf{a}_i^\top \text{Re}(\mathbf{R}) \mathbf{a}_i \right]^2 + \left[\sum_{i=1}^r \mathbf{a}_i^\top \text{Im}(\mathbf{R}) \mathbf{a}_i \right]^2 \right)^{1/2} \\ &\leq (\lambda_{\max}^2[\text{Re}(\mathbf{R})] + \lambda_{\max}^2[\text{Im}(\mathbf{R})])^{1/2} \text{tr}(\mathbf{A}^\top \mathbf{A}) \\ &\leq (2\lambda_{\max}[\text{Re}^2(\mathbf{R}) + \text{Im}^2(\mathbf{R})])^{1/2} \text{tr}(\mathbf{A}^\top \mathbf{A}) \\ &= 2^{1/2} \|\mathbf{R}\| |\text{tr}(\mathbf{A}^\top \mathbf{A})|. \quad \square \end{aligned}$$

Proof of Theorem 1. We prove part (i) first. Consider $\Psi_m^{(1)}(z) - \Psi^{(1)}(z) = I_1 + I_2$, where

$$I_1 = \Psi_m^{(1)}(z) - \frac{1}{p} \text{tr}(R_1(z)\Sigma_p), \quad I_2 = \frac{1}{p} \text{tr}(R_1(z)\Sigma_p) - \Psi^{(1)}(z),$$

with $R_1(z) = (\tilde{\Sigma}_1 - z\mathbf{I}_p)^{-1}$. A direct application of Lemma 2 of Ledoit and Péché (2011) implies that $p^{-1} \text{tr}(R_1(z)\Sigma_p)$ converges a.s. to $\Psi^{(1)}(z)$, and hence $I_2 \xrightarrow{a.s.} 0$. It remains to show that $I_1 \xrightarrow{a.s.} 0$ as well.

To this end, write $\mathbf{Y}_2 = (\mathbf{y}_{21}, \dots, \mathbf{y}_{2n_2})$, so that $\tilde{\Sigma}_2 = n_2^{-1} \sum_{k=1}^{n_2} \mathbf{y}_{2k} \mathbf{y}_{2k}^\top$. Also, write $\mathbf{y}_{2k} = \Sigma_p^{1/2} \mathbf{z}_{2k}$ for $k = 1, \dots, n_2$ with \mathbf{z}_{2k} being one of the \mathbf{z}_i 's in assumption (A1). Using the independence of the \mathbf{z}_{2k} 's and \mathbf{Y}_1 , we have for $k = 1, \dots, n_2$, $2 \leq q \leq 6$ and $z \in \mathbb{C}^+$,

$$\begin{aligned} & E \left\{ \left| \mathbf{z}_{2k}^\top \Sigma_p^{1/2} R_1(z) \Sigma_p^{1/2} \mathbf{z}_{2k} - \text{tr}(\Sigma_p^{1/2} R_1(z) \Sigma_p^{1/2}) \right|^q \middle| \mathbf{Y}_1 \right\} \\ & \leq K_q (E^{q/2} |z_{2k}|^4 \text{tr}^{q/2}(\Sigma_p R_1(z) \Sigma_p \overline{R_1(z)})) \\ & \quad + E |z_{2k}|^{2q} \text{tr}(\Sigma_p^{1/2} R_1(z) \Sigma_p \overline{R_1(z)} \Sigma_p^{1/2})^{q/2} \\ & \leq K_q (E^{q/2} |z_{2k}|^4 \|\Sigma_p\|^q p^{q/2} \text{Im}^{-q}(z)) \\ & \quad + E |z_{2k}|^{2q} p \lambda_{\max}^{q/2}(\Sigma_p^{1/2} R_1(z) \Sigma_p \overline{R_1(z)} \Sigma_p^{1/2}) \\ & \leq K_p \text{Im}^{-q}(z) p^{q/2} \|\Sigma_p\|^q (E^{q/2} |z_{2k}|^4 + p^{1-q/2} E |z_{2k}|^{2q}) = O(p^q), \end{aligned}$$

where the second line uses Lemma S.1, which is applicable by assumption (A1) on \mathbf{z}_{2k} , and the second last line is by Lemma S.3. The last line uses the assumptions in (A1). We have

$$\begin{aligned} E(|I_1|^6 | \mathbf{Y}_1) &= E \left(\left| \frac{1}{pn_2} \text{tr} \left(R_1(z) \sum_{k=1}^{n_2} \mathbf{y}_{2k} \mathbf{y}_{2k}^\top \right) - \frac{1}{p} \text{tr}(R_1(z)\Sigma_p) \right|^6 \middle| \mathbf{Y}_1 \right) \\ &= E \left(\left| \frac{1}{pn_2} \sum_{k=1}^{n_2} (\mathbf{z}_{2k}^\top \Sigma_p^{1/2} R_1(z) \Sigma_p^{1/2} \mathbf{z}_{2k} - \text{tr}(R_1(z)\Sigma_p)) \right|^6 \middle| \mathbf{Y}_1 \right). \end{aligned}$$

Define $g_k = \mathbf{z}_{2k}^\top \Sigma_p^{1/2} R_1(z) \Sigma_p^{1/2} \mathbf{z}_{2k} - \text{tr}(R_1(z)\Sigma_p)$. Then $E(g_i | \mathbf{Y}_1) = 0$, and $g_i | \mathbf{Y}_1$ is independent of $g_j | \mathbf{Y}_1$ for $i \neq j$. We then have the expansion

$$\begin{aligned} E(|I_1|^6) &= \frac{1}{p^6 n_2^6} \sum_{i_1, \dots, i_6=1}^{n_2} E(E(g_{i_1} \bar{g}_{i_2} g_{i_3} \bar{g}_{i_4} g_{i_5} \bar{g}_{i_6} | \mathbf{Y}_1)) \\ &\leq \frac{1}{p^6 n_2^6} E \left(\sum_{i=1}^{n_2} E(|g_i|^6 | \mathbf{Y}_1) + \sum_{i \neq j} [E(|g_i|^2 | \mathbf{Y}_1) E(|g_j|^4 | \mathbf{Y}_1) \right. \\ &\quad \left. + E(|g_i|^3 | \mathbf{Y}_1) E(|g_j|^3 | \mathbf{Y}_1)] \right. \\ &\quad \left. + \sum_{i \neq j \neq k} E(|g_i|^2 | \mathbf{Y}_1) E(|g_j|^2 | \mathbf{Y}_1) E(|g_k|^2 | \mathbf{Y}_1) \right) \quad (\text{S.1}) \\ &= O(p^{-6} n_2^{-6} \cdot [n_2 + 2n_2(n_2 - 1) + n_2(n_2 - 1)(n_2 - 2)] p^6) = O(n_2^{-3}), \end{aligned}$$

where the last line used $E(|g_i|^q | \mathbf{Y}_1) = O(p^q)$ proved before. Since $\sum_{n \geq 1} n_2^{-3} < \infty$ by assumption, we can apply the Borel-Cantelli lemma to conclude that $I_1 \xrightarrow{a.s.} 0$. This completes the proof of part (i).

Note also that since $I_1 \xrightarrow{a.s.} 0$, we have $\Psi_m^{(1)}(z) - p^{-1} \text{tr}(R_1(z) \Sigma_p) \xrightarrow{a.s.} 0$. Using the inversion formula (2.3), it is easy to show that the inverse Stieltjes transform for $\Psi_m^{(1)}(z)$ is $\Phi_m^{(1)}(x)$ in (2.13) on all points of continuity of $\Phi_m^{(1)}(x)$, and that for $p^{-1} \text{tr}(R_1(z) \Sigma_p)$ is $p^{-1} \sum_{i=1}^p \mathbf{p}_{1i}^\top \Sigma_p \mathbf{p}_{1i} \mathbf{1}_{\{\lambda_{1i} \leq x\}}$ on all points of continuity of the function. Hence, using equation (2.5) of Silverstein and Bai (1995), we conclude that

$$\Phi_m^{(1)}(x) - \frac{1}{p} \sum_{i=1}^p \mathbf{p}_{1i}^\top \Sigma_p \mathbf{p}_{1i} \mathbf{1}_{\{\lambda_{1i} \leq x\}} \xrightarrow{a.s.} 0$$

as $n_1, p \rightarrow \infty$ with $p/n_1 \rightarrow c_1 > 0$, which certainly include the case $c_1 = 1$. This proves the last part of the theorem.

The proof for part (ii) of the theorem can be found in the proof of Theorem 4 of Ledoit and P ech e (2011), from equations (45) to (50) by replacing γ there with $1/c_1$ and F with F_1 . Part (iii) is just an application of Theorem 4 of Ledoit and P ech e (2011). This completes the proof of the theorem. \square

Proof of Theorem 3. Let $\mathbf{Y}_2 = (\mathbf{y}_{21}, \dots, \mathbf{y}_{2n_2})$ with $\mathbf{y}_{2k} = \mathbf{A} \mathbf{x}_{2k} + \boldsymbol{\epsilon}_{2k}$, $k = 1, \dots, n_2$. Let $R_1(z) = (\tilde{\Sigma}_1 - z \mathbf{I}_p)^{-1}$. Then with $\Sigma_p = \mathbf{A} \Sigma_x \mathbf{A}^\top + \Sigma_\epsilon$ as in (3.18), we have

$$\begin{aligned} & \frac{1}{p} \text{tr}(R_1(z) \tilde{\Sigma}_2) - \frac{1}{p} \text{tr}(R_1(z) \Sigma_p) = D_1 + D_2 + D_3, \text{ where} \\ D_1 &= \frac{1}{pn_2} \sum_{k=1}^{n_2} \{ \mathbf{x}_{2k}^\top \mathbf{A}^\top R_1(z) \mathbf{A} \mathbf{x}_{2k} - \text{tr}(R_1(z) \mathbf{A} \Sigma_x \mathbf{A}^\top) \}, \\ D_2 &= \frac{2}{pn_2} \sum_{k=1}^{n_2} \mathbf{x}_{2k}^\top \mathbf{A}^\top R_1(z) \boldsymbol{\epsilon}_{2k}, \\ D_3 &= \frac{1}{pn_2} \sum_{k=1}^{n_2} \{ \boldsymbol{\epsilon}_{2k}^\top R_1(z) \boldsymbol{\epsilon}_{2k} - \text{tr}(R_1(z) \Sigma_\epsilon) \}. \end{aligned}$$

Consider D_1 first. Using assumption (F1), and defining

$$g_i = \mathbf{x}_{2k}^{*\top} \Sigma_x^{1/2} \mathbf{A}^\top R_1(z) \mathbf{A} \Sigma_x^{1/2} \mathbf{x}_{2k}^* - \text{tr}(R_1(z) \mathbf{A} \Sigma_x \mathbf{A}^\top),$$

for $2 \leq q \leq 6$ and $z \in \mathbb{C}^+$, we have by Lemma S.1,

$$\begin{aligned} E(|g_k|^q | \mathbf{Y}_1) &\leq K_q (E^{q/2} |x_{2k}^*|^4 \text{tr}^{q/2}(\Sigma_x^{1/2} \mathbf{A}^\top R_1(z) \mathbf{A} \Sigma_x \mathbf{A}^\top \overline{R_1(z)} \mathbf{A} \Sigma_x^{1/2}) \\ &\quad + E |x_{2k}^*|^{2q} \text{tr}(\Sigma_x^{1/2} \mathbf{A}^\top R_1(z) \mathbf{A} \Sigma_x \mathbf{A}^\top \overline{R_1(z)} \mathbf{A} \Sigma_x^{1/2})^{q/2}) \\ &\leq K_q \left\| \Sigma_x^{1/2} \mathbf{A}^\top R_1(z) \mathbf{A} \Sigma_x \mathbf{A}^\top \overline{R_1(z)} \mathbf{A} \Sigma_x^{1/2} \right\|^{q/2} \\ &\quad \cdot (E^{q/2} |x_{2k}^*|^{4r^{q/2}} + r E |x_{2k}^*|^{2q}) \\ &\leq K_q \left\| \Sigma_x \right\|^q \left\| \mathbf{A} \right\|^{2q} \text{Im}^{-q}(z) (E^{q/2} |x_{2k}^*|^{4r^{q/2}} + r E |x_{2k}^*|^{2q}) = O(p^q), \end{aligned}$$

where we have used $\|\mathbf{A}\| = O(p^{1/2})$ implied in assumption (F2). Using the fact that $E(g_k|\mathbf{Y}_1) = 0$ and $g_i|\mathbf{Y}_1$ is independent of $g_j|\mathbf{Y}_1$ for $i \neq j$ by assumption (F1), similar to the expansion in (S.1), we have

$$E(|D_1|^6) = \frac{1}{p^6 n_2^6} \sum_{i_1, \dots, i_6=1}^{n_2} E(E(g_{i_1} \bar{g}_{i_2} g_{i_3} \bar{g}_{i_4} g_{i_5} \bar{g}_{i_6} | \mathbf{Y}_1)) = O(n_2^{-3}),$$

so that $D_1 \xrightarrow{a.s.} 0$ by the Borel-Cantelli lemma, as $\sum_{n \geq 1} n_2^{-3} < \infty$ by assumption.

Now consider D_2 . Define $g_k = \mathbf{x}_{2k}^T \mathbf{A}^T R_1(z) \boldsymbol{\epsilon}_{2k}$. We have for $0 < q \leq 12$ and $z \in \mathbb{C}^+$,

$$E(|g_k|^q | \mathbf{Y}_1) \leq E\|\mathbf{x}_{2k}^*\|^q \|\boldsymbol{\Sigma}_x\|^{q/2} \|\mathbf{A}\|^q E\|\boldsymbol{\xi}_{2k}\|^q \|\boldsymbol{\Sigma}_\epsilon\|^{q/2} \text{Im}^{-q}(z) = O(p^q),$$

by assumptions (F1) and (F2), and lemma S.2. Hence, using the fact that $E(g_k|\mathbf{Y}_1) = 0$ and $g_i|\mathbf{Y}_1$ is independent of $g_j|\mathbf{Y}_1$ for $i \neq j$, an expansion similar to (S.1) leads again to

$$E(|D_2|^6) = O(n_2^{-3}),$$

so that by the Borel-Cantelli lemma, we have $D_2 \xrightarrow{a.s.} 0$.

Finally, if we define $g_k = \boldsymbol{\xi}_{2k}^T \boldsymbol{\Sigma}_\epsilon^{1/2} R_1(z) \boldsymbol{\Sigma}_\epsilon^{1/2} \boldsymbol{\xi}_{2k} - \text{tr}(R_1(z) \boldsymbol{\Sigma}_\epsilon)$, then following exactly the same arguments as in the proof of Theorem 1, we can conclude

$$E(|D_3|^6) = O(n_2^{-3}),$$

so that through the Borel-Cantelli lemma we have $D_3 \xrightarrow{a.s.} 0$.

Hence, we have proved

$$\frac{1}{p} \text{tr}[R_1(z) \tilde{\boldsymbol{\Sigma}}_2] - \frac{1}{p} \text{tr}[R_1(z) \boldsymbol{\Sigma}_p] \xrightarrow{a.s.} 0.$$

Then, using equation (2.5) of Silverstein and Bai (1995), we conclude the almost sure convergence of the inverse Stieltjes transform on all points of continuity x , that is,

$$\frac{1}{p} \sum_{i=1}^p \mathbf{p}_{1i}^T \tilde{\boldsymbol{\Sigma}}_2 \mathbf{p}_{1i} \mathbf{1}_{\{\lambda_{1i} \leq x\}} - \frac{1}{p} \sum_{i=1}^p \mathbf{p}_{1i}^T \boldsymbol{\Sigma}_p \mathbf{p}_{1i} \mathbf{1}_{\{\lambda_{1i} \leq x\}} \xrightarrow{a.s.} 0. \quad \square$$

Proof of Lemma 1. We first assume (A1)' and (A2)'. For $j = 1, \dots, n_2$ and $i = 1, \dots, p$, define

$$g_{ij} = \frac{\mathbf{z}_{2j}^T \boldsymbol{\Sigma}_p^{1/2} \mathbf{p}_{1i} \mathbf{p}_{1i}^T \boldsymbol{\Sigma}_p^{1/2} \mathbf{z}_{2j} - \text{tr}(\boldsymbol{\Sigma}_p^{1/2} \mathbf{p}_{1i} \mathbf{p}_{1i}^T \boldsymbol{\Sigma}_p^{1/2})}{\mathbf{p}_{1i}^T \boldsymbol{\Sigma}_p \mathbf{p}_{1i}}.$$

Then for $2 \leq q \leq 10$, using Lemma S.1 and assumption (A1)',

$$\begin{aligned} E(|g_{ij}|^q | \mathbf{Y}_1) &\leq (\mathbf{p}_{1i}^T \boldsymbol{\Sigma}_p \mathbf{p}_{1i})^{-q} K_q \left(E^{q/2} |z_{2j}|^4 \text{tr}^q(\boldsymbol{\Sigma}_p^{1/2} \mathbf{p}_{1i} \mathbf{p}_{1i}^T \boldsymbol{\Sigma}_p^{1/2}) \right. \\ &\quad \left. + E |z_{2j}|^{2q} \text{tr}(\boldsymbol{\Sigma}_p^{1/2} \mathbf{p}_{1i} \mathbf{p}_{1i}^T \boldsymbol{\Sigma}_p^{1/2})^q \right) \\ &= K_q (E^{q/2} |z_{2j}|^4 + E |z_{2j}|^{2q}) = O(1). \end{aligned}$$

At the same time, $E(g_{ij}|\mathbf{Y}_1) = 0$ and $g_{i_1 j_1}|\mathbf{Y}_1$ is independent of $g_{i_2 j_2}|\mathbf{Y}_1$ whenever $j_1 \neq j_2$. Hence,

$$\begin{aligned}
& E \left\{ \max_{1 \leq i \leq p} \left| \frac{\mathbf{p}_{1i}^\top \tilde{\Sigma}_2 \mathbf{p}_{1i} - \mathbf{p}_{1i}^\top \Sigma_p \mathbf{p}_{1i}}{\mathbf{p}_{1i}^\top \Sigma_p \mathbf{p}_{1i}} \right| \right\}^{10} = E \left\{ \max_{1 \leq i \leq p} \left| \frac{1}{n_2} \sum_{j=1}^{n_2} g_{ij} \right|^{10} \right\} \\
& \leq \sum_{i=1}^p \frac{1}{n_2^{10}} \sum_{j_1 \cdots j_{10}=1}^{n_2} E(E\{g_{ij_1} \cdots g_{ij_{10}}|\mathbf{Y}_1\}) \\
& \leq E \left\{ \sum_{i=1}^p \frac{1}{n_2^{10}} \left(\sum_{j=1}^{n_2} E(g_{ij}^{10}|\mathbf{Y}_1) + \sum_{\substack{q_1, q_2 \neq 1 \\ q_1 + q_2 = 10}} \sum_{j_1 \neq j_2} E(g_{ij_1}^{q_1}|\mathbf{Y}_1) E(g_{ij_2}^{q_2}|\mathbf{Y}_1) \right. \right. \\
& + \sum_{\substack{q_1, q_2, q_3 \neq 1 \\ q_1 + q_2 + q_3 = 10}} \sum_{j_1 \neq j_2 \neq j_3} \prod_{k=1}^3 E(g_{ij_k}^{q_k}|\mathbf{Y}_1) + \sum_{\substack{q_1, \dots, q_4 \neq 1 \\ q_1 + \dots + q_4 = 10}} \sum_{j_1 \neq j_2 \neq j_3 \neq j_4} \prod_{k=1}^4 E(g_{ij_k}^{q_k}|\mathbf{Y}_1) \\
& \left. \left. + \sum_{j_1 \neq \dots \neq j_5} \prod_{k=1}^5 E(g_{ij_k}^2|\mathbf{Y}_1) \right) \right\} \\
& = O(p n_2^{-10} [n_2 + 4P_2^{n_2} + 4P_3^{n_2} + 2P_4^{n_2} + P_5^{n_2}]) \\
& = O(p n_2^{-5}), \tag{S.2}
\end{aligned}$$

where $P_r^{n_2} = n_2(n_2 - 1) \cdots (n_2 - r + 1)$. Since $\sum_{n \geq 1} p n_2^{-5} < \infty$ by assumption, the Borel-Cantelli lemma completes the proof of this part.

Now consider data from a factor model, with assumption (F1)' satisfied. We can decompose

$$\begin{aligned}
& \frac{\mathbf{p}_{1i}^\top \tilde{\Sigma}_2 \mathbf{p}_{1i} - \mathbf{p}_{1i}^\top \Sigma_p \mathbf{p}_{1i}}{\mathbf{p}_{1i}^\top \Sigma_p \mathbf{p}_{1i}} = D_1 + D_2 + D_3, \quad \text{where} \\
& D_1 = \frac{1}{n_2} \sum_{j=1}^{n_2} g_{ij}, \quad g_{ij} = \frac{\mathbf{x}_{2j}^{*\top} \Sigma_x^{1/2} \mathbf{A}^\top \mathbf{p}_{1i} \mathbf{p}_{1i}^\top \mathbf{A} \Sigma_x^{1/2} \mathbf{x}_{2j}^* - \mathbf{p}_{1i}^\top \mathbf{A} \Sigma_x \mathbf{A}^\top \mathbf{p}_{1i}}{\mathbf{p}_{1i}^\top \Sigma_p \mathbf{p}_{1i}}, \\
& D_2 = \frac{1}{n_2} \sum_{j=1}^{n_2} d_{ij}, \quad d_{ij} = \frac{\boldsymbol{\xi}_{2j}^\top \Sigma_\epsilon^{1/2} \mathbf{p}_{1i} \mathbf{p}_{1i}^\top \Sigma_\epsilon^{1/2} \boldsymbol{\xi}_{2j} - \mathbf{p}_{1i}^\top \Sigma_\epsilon \mathbf{p}_{1i}}{\mathbf{p}_{1i}^\top \Sigma_p \mathbf{p}_{1i}}, \\
& D_3 = \frac{2}{n_2} \sum_{j=1}^{n_2} h_{ij}, \quad h_{ij} = \frac{\mathbf{x}_{2j}^{*\top} \Sigma_x^{1/2} \mathbf{A}^\top \mathbf{p}_{1i} \mathbf{p}_{1i}^\top \Sigma_\epsilon^{1/2} \boldsymbol{\xi}_{2j}}{\mathbf{p}_{1i}^\top \Sigma_p \mathbf{p}_{1i}}.
\end{aligned}$$

Using Lemma S.1, for $2 \leq q \leq 10$,

$$\begin{aligned}
& E(|g_{ij}|^q|\mathbf{Y}_1) \leq (\mathbf{p}_{1i}^\top \Sigma_p \mathbf{p}_{1i})^{-q} K_q (\mathbf{p}_{1i}^\top \mathbf{A} \Sigma_x \mathbf{A}^\top \mathbf{p}_{1i})^q (E^{q/2} |x_{2j}^*|^4 + E |x_{2j}^*|^{2q}) \\
& \leq K_q (E^{q/2} |x_{2j}^*|^4 + E |x_{2j}^*|^{2q}) = O(1); \\
& E(|d_{ij}|^q|\mathbf{Y}_1) \leq (\mathbf{p}_{1i}^\top \Sigma_p \mathbf{p}_{1i})^{-q} K_q (\mathbf{p}_{1i}^\top \Sigma_\epsilon \mathbf{p}_{1i})^q (E^{q/2} |\xi_{2j}|^4 + E |\xi_{2j}|^{2q}) \\
& \leq K_q (E^{q/2} |\xi_{2j}|^4 + E |\xi_{2j}|^{2q}) = O(1).
\end{aligned}$$

For $4 \leq q \leq 10$, using the c_r -inequality,

$$E(|h_{ij}|^q | \mathbf{Y}_1) = E \left(\left(g_{ij} + \frac{\mathbf{p}_{1i}^\top \mathbf{A} \boldsymbol{\Sigma}_x \mathbf{A}^\top \mathbf{p}_{1i}}{\mathbf{p}_{1i}^\top \boldsymbol{\Sigma}_p \mathbf{p}_{1i}} \right)^{q/2} \middle| \mathbf{Y}_1 \right) \\ \cdot E \left(\left(d_{ij} + \frac{\mathbf{p}_{1i}^\top \boldsymbol{\Sigma}_\epsilon \mathbf{p}_{1i}}{\mathbf{p}_{1i}^\top \boldsymbol{\Sigma}_p \mathbf{p}_{1i}} \right)^{q/2} \middle| \mathbf{Y}_1 \right), \text{ with}$$

$$E \left(\left(g_{ij} + \frac{\mathbf{p}_{1i}^\top \mathbf{A} \boldsymbol{\Sigma}_x \mathbf{A}^\top \mathbf{p}_{1i}}{\mathbf{p}_{1i}^\top \boldsymbol{\Sigma}_p \mathbf{p}_{1i}} \right)^{q/2} \middle| \mathbf{Y}_1 \right) \leq C_{q/2} (E(|g_{ij}|^{q/2} | \mathbf{Y}_1) + 1) = O(1), \\ E \left(\left(\frac{\boldsymbol{\xi}_{2j}^\top \boldsymbol{\Sigma}_\epsilon^{1/2} \mathbf{p}_{1i} \mathbf{p}_{1i}^\top \boldsymbol{\Sigma}_\epsilon^{1/2} \boldsymbol{\xi}_{2j}}{\mathbf{p}_{1i}^\top \boldsymbol{\Sigma}_p \mathbf{p}_{1i}} \right)^{q/2} \middle| \mathbf{Y}_1 \right) \leq C_{q/2} (E(|d_{ij}|^{q/2} | \mathbf{Y}_1) + 1) = O(1),$$

where $C_{q/2}$ is a constant. Hence, $E(|h_{ij}|^q | \mathbf{Y}_1) = O(1)$ for $q \geq 4$. For $2 \leq q < 4$, it is easy to see that $E(|h_{ij}|^q | \mathbf{Y}_1) \leq E^{1/2}(|h_{ij}|^{2q} | \mathbf{Y}_1) = O(1)$, and hence $E(|h_{ij}|^q | \mathbf{Y}_1) = O(1)$ for $2 \leq q \leq 10$. Then since $E(g_{ij} | \mathbf{Y}_1) = 0$ and $g_{i_1 j_1} | \mathbf{Y}_1$ is independent of $g_{i_2 j_2} | \mathbf{Y}_1$ whenever $j_1 \neq j_2$, an expansion exactly like (S.2) will show that

$$E \left(\left| \max_{1 \leq i \leq p} |D_1|^{10} \right| \right) = O(pn_2^{-5}),$$

showing that $\max_{1 \leq i \leq p} |D_1| \xrightarrow{a.s.} 0$, since $\sum_{n \geq 1} pn_2^{-5} < \infty$ by assumption. Using the same expansion (S.2) and the arguments above, we can also show that $\max_{1 \leq i \leq p} |D_2|, \max_{1 \leq i \leq p} |D_3| \xrightarrow{a.s.} 0$. This completes the proof of the lemma. \square

Before proving Theorem 5, we state and prove a lemma first.

Lemma S.4 *Let assumptions (A1) to (A4) be satisfied. Then assuming $\boldsymbol{\Sigma}_p \neq \sigma^2 \mathbf{I}$,*

$$\frac{\|\widehat{\boldsymbol{\Sigma}}_{\text{Ideal}, m} - \boldsymbol{\Sigma}_p\|_F^2}{\|\widehat{\boldsymbol{\Sigma}}_{\text{Ideal}} - \boldsymbol{\Sigma}_p\|_F^2} \xrightarrow{a.s.} 1, \quad \frac{\log \det(\widehat{\boldsymbol{\Sigma}}_{\text{Ideal}, m} \boldsymbol{\Sigma}_p^{-1})}{\log \det(\widehat{\boldsymbol{\Sigma}}_{\text{Ideal}} \boldsymbol{\Sigma}_p^{-1})} \xrightarrow{a.s.} 1.$$

Proof of Lemma S.4. Consider the Frobenius loss first. For λ_i the i th largest eigenvalue of \mathbf{S}_n with corresponding eigenvector \mathbf{p}_i ,

$$p^{-1} \|\widehat{\boldsymbol{\Sigma}}_{\text{Ideal}} - \boldsymbol{\Sigma}_p\|_F^2 = p^{-1} \text{tr}(\text{diag}(\mathbf{P}^\top \boldsymbol{\Sigma}_p \mathbf{P}) - \mathbf{P}^\top \boldsymbol{\Sigma}_p \mathbf{P})^2 \\ = p^{-1} \sum_{i=1}^n (\mathbf{p}_i^\top \boldsymbol{\Sigma}_p \mathbf{p}_i)^2 - 2p^{-1} \text{tr}(\mathbf{P}^\top \boldsymbol{\Sigma}_p \mathbf{P} \text{diag}(\mathbf{P}^\top \boldsymbol{\Sigma}_p \mathbf{P})) \\ + p^{-1} \text{tr}(\boldsymbol{\Sigma}_p^2) \\ = p^{-1} \text{tr}(\boldsymbol{\Sigma}_p^2) - p^{-1} \sum_{i=1}^p (\mathbf{p}_i^\top \boldsymbol{\Sigma}_p \mathbf{p}_i)^2,$$

so that similarly, $p^{-1}\|\widehat{\Sigma}_{\text{Ideal},m} - \Sigma_p\|_F^2 = p^{-1}\text{tr}(\Sigma_p^2) - p^{-1}\sum_{i=1}^p(\mathbf{P}_{1i}^T \Sigma_p \mathbf{P}_{1i})^2$. But with assumption (A3), we have

$$p^{-1}\text{tr}(\Sigma_p^2) = p^{-1}\sum_{i=1}^p \tau_{n,i}^2 \xrightarrow{a.s.} \int \tau^2 dH(\tau).$$

Moreover, with the convergence result in Theorem 4 of Ledoit and P ech e (2011) indicating that (see equation (2.7) as well) for all $x \in \mathbb{R}$, we have

$$\Delta_p(x) = p^{-1}\sum_{i=1}^p \mathbf{P}_i^T \Sigma_p \mathbf{P}_i \mathbf{1}_{\{\lambda_i \leq x\}} \xrightarrow{a.s.} \Delta(x) = \int_{-\infty}^x \delta(\lambda) dF(\lambda),$$

where $F(\cdot)$ is such that

$$F_p(\lambda) = p^{-1}\sum_{i=1}^p \mathbf{1}_{\{\lambda_i \leq \lambda\}} \xrightarrow{a.s.} F(\lambda),$$

we have for any continuous function $g(\cdot)$ over the positive real line,

$$p^{-1}\sum_{i=1}^p g(\mathbf{P}_i^T \Sigma_p \mathbf{P}_i) \mathbf{1}_{\{\lambda_i \leq x\}} \xrightarrow{a.s.} \int_{-\infty}^x g(\delta(\lambda)) dF(\lambda). \quad (\text{S.3})$$

With this, it is easy to see that then

$$p^{-1}\|\widehat{\Sigma}_{\text{Ideal}} - \Sigma_p\|_F^2 \xrightarrow{a.s.} \int \tau^2 dH(\tau) - \int \delta^2(\lambda) dF(\lambda),$$

which is nonzero in general, except when $\Sigma_p = \sigma^2 \mathbf{I}$. At the same time, we have (see equation (2.9) and the descriptions therein as well)

$$p^{-1}\sum_{i=1}^p \mathbf{P}_{1i}^T \Sigma_p \mathbf{P}_{1i} \mathbf{1}_{\{\lambda_{1i} \leq x\}} \xrightarrow{a.s.} \int_{-\infty}^x \delta_1(\lambda) dF_1(\lambda),$$

where $F_1(\cdot)$ is such that

$$F_{1p}(\lambda) = p^{-1}\sum_{i=1}^p \mathbf{1}_{\{\lambda_{1i} \leq \lambda\}} \xrightarrow{a.s.} F_1(\lambda).$$

But since $m/n = n_1/n \rightarrow 1$, we have $p/n_1, p/n$ both going to the same limit $c_1 = c > 0$. Theorem 4.1 of Bai and Silverstein (2010) tells us then both F_p and F_{1p} converges to the same limit almost surely under assumptions (A1) to (A4). Hence $F = F_1$ almost surely, implying $\delta_1(\cdot) = \delta(\cdot)$ almost surely (See Remark 1 in section 2.2 as well). This immediately implies that, similar to (S.3),

$$p^{-1}\sum_{i=1}^p (\mathbf{P}_{1i}^T \Sigma_p \mathbf{P}_{1i})^2 \mathbf{1}_{\{\lambda_{1i} \leq x\}} \xrightarrow{a.s.} \int_{-\infty}^x \delta_1^2(\lambda) dF_1(\lambda) = \int_{-\infty}^x \delta^2(\lambda) dF(\lambda),$$

and hence

$$\begin{aligned} p^{-1} \|\widehat{\Sigma}_{\text{Ideal},m} - \Sigma_p\|_F^2 &\xrightarrow{a.s.} \int \tau^2 dH(\tau) - \int \delta_1^2(\lambda) dF_1(\lambda) \\ &= \int \tau^2 dH(\tau) - \int \delta^2(\lambda) dF(\lambda). \end{aligned}$$

We immediately have

$$\frac{\|\widehat{\Sigma}_{\text{Ideal},m} - \Sigma_p\|_F^2}{\|\widehat{\Sigma}_{\text{Ideal}} - \Sigma_p\|_F^2} \xrightarrow{a.s.} 1,$$

which completes the proof for the Frobenius loss.

For the remaining part, note that it is not difficult to show the inverse Stein's loss

$$\begin{aligned} p^{-1} SL(\Sigma_p, \widehat{\Sigma}_{\text{Ideal}}) &= p^{-1} \log \det(\widehat{\Sigma}_{\text{Ideal}} \Sigma_p^{-1}) \\ &= p^{-1} \sum_{i=1}^p \log(\mathbf{p}_i^\top \Sigma_p \mathbf{p}_i) - p^{-1} \sum_{i=1}^p \log(\tau_{n,i}). \end{aligned}$$

Similar arguments as before show that

$$\begin{aligned} p^{-1} \sum_{i=1}^p \log(\mathbf{p}_i^\top \Sigma_p \mathbf{p}_i) &\xrightarrow{a.s.} \int \log(\delta(\lambda)) dF(\lambda), \\ p^{-1} \sum_{i=1}^p \log(\tau_{n,i}) &\xrightarrow{a.s.} \int \log(\tau) dH(\tau). \end{aligned}$$

Hence

$$\begin{aligned} p^{-1} SL(\Sigma_p, \widehat{\Sigma}_{\text{Ideal}}) &= p^{-1} \log \det(\widehat{\Sigma}_{\text{Ideal}} \Sigma_p^{-1}) \\ &\xrightarrow{a.s.} \int \log(\delta(\lambda)) dF(\lambda) - \int \log(\tau) dH(\tau), \end{aligned}$$

which is nonzero in general except when $\Sigma_p = \sigma^2 \mathbf{I}$. Moreover, similar arguments as in the case for the Frobenius loss show that

$$\begin{aligned} p^{-1} \sum_{i=1}^p \log(\mathbf{p}_{1i}^\top \Sigma_p \mathbf{p}_{1i}) \mathbf{1}_{\{\lambda_{1i} \leq x\}} &\xrightarrow{a.s.} \int_{-\infty}^x \log(\delta_1(\lambda)) dF_1(\lambda) \\ &= \int_{-\infty}^x \log(\delta(\lambda)) dF(\lambda). \end{aligned}$$

Hence,

$$\begin{aligned} p^{-1} SL(\Sigma_p, \widehat{\Sigma}_{\text{Ideal},m}) &= p^{-1} \log \det(\widehat{\Sigma}_{\text{Ideal},m} \Sigma_p^{-1}) \\ &\xrightarrow{a.s.} \int \log(\delta_1(\lambda)) dF_1(\lambda) - \int \log(\tau) dH(\tau) \\ &= \int \log(\delta(\lambda)) dF(\lambda) - \int \log(\tau) dH(\tau). \end{aligned}$$

This completes the proof of the lemma. \square

Proof of Theorem 5. We first assume (A1)', (A2)', (A3) and (A4). Since we have

$$\begin{aligned} p^{-1} \|\widehat{\Sigma}_m - \Sigma_p\|_F^2 &= p^{-1} \|\text{diag}(\mathbf{P}_1^T \widetilde{\Sigma}_2 \mathbf{P}_1) - \mathbf{P}_1^T \Sigma_p \mathbf{P}_1\|_F^2 \\ &= p^{-1} \sum_{i=1}^p (\mathbf{p}_{1i}^T \widetilde{\Sigma}_2 \mathbf{p}_{1i} - \mathbf{p}_{1i}^T \Sigma_p \mathbf{p}_{1i})^2 + p^{-1} \|\widehat{\Sigma}_{\text{Ideal},m} - \Sigma_p\|_F^2, \end{aligned}$$

we can express the efficiency loss in (4.2) with respect to the Frobenius loss as

$$\begin{aligned} EL(\Sigma_p, \widehat{\Sigma}_m) &= 1 - \left(\frac{p^{-1} \sum_{i=1}^p (\mathbf{p}_{1i}^T \widetilde{\Sigma}_2 \mathbf{p}_{1i} - \mathbf{p}_{1i}^T \Sigma_p \mathbf{p}_{1i})^2}{p^{-1} \|\widehat{\Sigma}_{\text{Ideal},m} - \Sigma_p\|_F^2} \right. \\ &\quad \left. + \frac{p^{-1} \|\widehat{\Sigma}_{\text{Ideal},m} - \Sigma_p\|_F^2}{p^{-1} \|\widehat{\Sigma}_{\text{Ideal},m} - \Sigma_p\|_F^2} \right)^{-1}. \end{aligned}$$

With $m/n \rightarrow 1$, we have $p/m \rightarrow c > 0$, same as p/n does. The conclusion from lemma S.4 that the ratio of loss for the ideal estimators being almost surely 1 then implies that $EL(\widehat{\Sigma}_m - \Sigma_p) \xrightarrow{a.s.} 0$, if we can show that $p^{-1} \sum_{i=1}^p (\mathbf{p}_{1i}^T \widetilde{\Sigma}_2 \mathbf{p}_{1i} - \mathbf{p}_{1i}^T \Sigma_p \mathbf{p}_{1i})^2 \xrightarrow{a.s.} 0$. But together with assumption $\sum_{n \geq 1} pn_2^{-5} < \infty$, we can apply Lemma 1, so that

$$\begin{aligned} p^{-1} \sum_{i=1}^p (\mathbf{p}_{1i}^T \widetilde{\Sigma}_2 \mathbf{p}_{1i} - \mathbf{p}_{1i}^T \Sigma_p \mathbf{p}_{1i})^2 &\leq \left(\max_{1 \leq i \leq p} \left| \frac{\mathbf{p}_{1i}^T \widetilde{\Sigma}_2 \mathbf{p}_{1i} - \mathbf{p}_{1i}^T \Sigma_p \mathbf{p}_{1i}}{\mathbf{p}_{1i}^T \Sigma_p \mathbf{p}_{1i}} \right| \right)^2 \\ &\quad \cdot \max_{1 \leq i \leq p} (\mathbf{p}_{1i}^T \Sigma_p \mathbf{p}_{1i})^2 \xrightarrow{a.s.} 0, \end{aligned}$$

since $\max_{1 \leq i \leq p} (\mathbf{p}_{1i}^T \Sigma_p \mathbf{p}_{1i})^2 \leq \|\Sigma_p\|^2 = O(1)$ by assumption (A2)'. Hence, $EL(\Sigma_p, \widehat{\Sigma}_m) \xrightarrow{a.s.} 0$ in this case.

Now consider the inverse Stein's loss. We have

$$\begin{aligned} SL(\Sigma_p, \widehat{\Sigma}_m) &= \text{tr}(\mathbf{P}_1^T \Sigma_p \mathbf{P}_1 \text{diag}^{-1}(\mathbf{P}_1^T \widetilde{\Sigma}_2 \mathbf{P}_1)) \\ &\quad - \log \det(\mathbf{P}_1^T \Sigma_p \mathbf{P}_1 \text{diag}^{-1}(\mathbf{P}_1^T \widetilde{\Sigma}_2 \mathbf{P}_1)) - p \\ &= \sum_{i=1}^p \frac{\mathbf{p}_{1i}^T \Sigma_p \mathbf{p}_{1i}}{\mathbf{p}_{1i}^T \widetilde{\Sigma}_2 \mathbf{p}_{1i}} - p + \sum_{i=1}^p \log \left(\frac{\mathbf{p}_{1i}^T \widetilde{\Sigma}_2 \mathbf{p}_{1i}}{\mathbf{p}_{1i}^T \Sigma_p \mathbf{p}_{1i}} \right) + SL(\Sigma_p, \widehat{\Sigma}_{\text{Ideal},m}) \\ &= \sum_{i=1}^p h_1(g_i) + \sum_{i=1}^p h_2(g_i) + SL(\Sigma_p, \widehat{\Sigma}_{\text{Ideal},m}), \end{aligned}$$

where we define g_i for $i = 1, \dots, p$, and the functions h_1 and h_2 , which are differentiable on $(-1, \infty)$, to be

$$g_i = \frac{\mathbf{p}_{1i}^T \widetilde{\Sigma}_2 \mathbf{p}_{1i} - \mathbf{p}_{1i}^T \Sigma_p \mathbf{p}_{1i}}{\mathbf{p}_{1i}^T \Sigma_p \mathbf{p}_{1i}}, \quad h_1(x) = \frac{1}{1+x} - 1, \quad h_2(x) = \log(1+x).$$

With $SL(\Sigma_p, \widehat{\Sigma}_{\text{Ideal}, m})/SL(\Sigma_p, \widehat{\Sigma}_{\text{Ideal}}) \xrightarrow{a.s.} 1$ by lemma S.4,

$$EL(\Sigma_p, \widehat{\Sigma}_m) = 1 - \left(\frac{p^{-1} \sum_{i=1}^p h_1(g_i)}{p^{-1} SL(\Sigma_p, \widehat{\Sigma}_{\text{Ideal}})} + \frac{p^{-1} \sum_{i=1}^p h_2(g_i)}{p^{-1} SL(\Sigma_p, \widehat{\Sigma}_{\text{Ideal}})} + \frac{p^{-1} SL(\Sigma_p, \widehat{\Sigma}_{\text{Ideal}, m})}{p^{-1} SL(\Sigma_p, \widehat{\Sigma}_{\text{Ideal}})} \right)^{-1} \xrightarrow{a.s.} 0,$$

if we can show further that $p^{-1} \sum_{i=1}^p h_j(g_i) \xrightarrow{a.s.} 0$ for $j = 1, 2$. To this end, by Lemma 1, g_i is almost surely 0 for all i . Hence, in particular, for n large enough, we have $|g_i| \leq B < 1$ almost surely, so that $h_j(g_i)$ is well defined for all i . This implies that for $j = 1, 2$,

$$p^{-1} \sum_{i=1}^p h_j(g_i) \leq \max_{1 \leq i \leq p} h_j(g_i) = \max_{1 \leq i \leq p} |g_i| |h'_j(\eta_i)| \xrightarrow{a.s.} 0,$$

where η_i lies almost surely in $[-B, B]$ for all i for large enough n , so that the $h'_j(\eta_i)$'s are bounded almost surely for all i . This completes the proof of the theorem. \square

Proof of Theorem 6. For the Frobenius loss,

$$\begin{aligned} \|\widehat{\Sigma}_{m, M} - \Sigma_p\|_F^2 &= \left\| \frac{1}{M} \sum_{i=1}^M (\widehat{\Sigma}_m^{(i)} - \Sigma_p) \right\|_F^2 \leq \left(\frac{1}{M} \sum_{i=1}^M \|\widehat{\Sigma}_m^{(i)} - \Sigma_p\|_F \right)^2 \\ &\leq \frac{1}{M} \sum_{i=1}^M \|\widehat{\Sigma}_m^{(i)} - \Sigma_p\|_F^2, \end{aligned}$$

so that

$$EL(\Sigma_p, \widehat{\Sigma}_{m, M}) \leq 1 - \frac{\|\widehat{\Sigma}_{\text{Ideal}} - \Sigma_p\|_F^2}{\frac{1}{M} \sum_{i=1}^M \|\widehat{\Sigma}_m^{(i)} - \Sigma_p\|_F^2} = 1 - \frac{1}{\frac{1}{M} \sum_{i=1}^M \frac{1}{1 - EL(\Sigma_p, \widehat{\Sigma}_m^{(i)})}} \xrightarrow{a.s.} 0,$$

since $EL(\Sigma_p, \widehat{\Sigma}_m^{(i)}) \xrightarrow{a.s.} 0$ by Theorem 5. This completes the proof for the Frobenius loss.

For the inverse Stein's loss, note that for any estimator $\widehat{\Sigma}$,

$$SL(\Sigma_p, \widehat{\Sigma}) = \text{tr} \left((\Sigma_p - \widehat{\Sigma}) \widehat{\Sigma}^{-1} \right) + \log \det(\widehat{\Sigma} \Sigma_p^{-1}).$$

Hence

$$\begin{aligned}
SL(\Sigma_p, \widehat{\Sigma}_{m,M}) &= \frac{1}{M} \sum_{i=1}^M \text{tr} \left((\Sigma_p - \widehat{\Sigma}_m^{(i)}) \widehat{\Sigma}_{m,M}^{-1} \right) + \log \det \left(\frac{1}{M} \sum_{i=1}^M \widehat{\Sigma}_m^{(i)} \Sigma_p^{-1} \right) \\
&\leq \frac{1}{M} \sum_{i=1}^M \text{tr} \left((\Sigma_p - \widehat{\Sigma}_m^{(i)}) \widehat{\Sigma}_{m,M}^{-1} \right) + \frac{1}{M} \sum_{i=1}^M \log \det (\widehat{\Sigma}_m^{(i)} \Sigma_p^{-1}) \\
&= \frac{1}{M} \sum_{i=1}^M \text{tr} \left((\Sigma_p - \widehat{\Sigma}_m^{(i)}) (\widehat{\Sigma}_{m,M}^{-1} - (\widehat{\Sigma}_m^{(i)})^{-1}) \right) + \frac{1}{M} \sum_{i=1}^M SL(\Sigma_p, \widehat{\Sigma}_m^{(i)}),
\end{aligned}$$

where the second line comes from the convexity of $\log \det(\mathbf{X})$ where \mathbf{X} is positive semi-definite. Then

$$\begin{aligned}
EL(\Sigma_p, \widehat{\Sigma}_{m,M}) &\leq 1 - \left(\frac{\frac{1}{pM} \sum_{i=1}^M \text{tr} \left((\Sigma_p - \widehat{\Sigma}_m^{(i)}) (\widehat{\Sigma}_{m,M}^{-1} - (\widehat{\Sigma}_m^{(i)})^{-1}) \right)}{p^{-1} SL(\Sigma_p, \widehat{\Sigma}_{\text{Ideal}})} \right. \\
&\quad \left. + \frac{1}{M} \sum_{i=1}^M \frac{1}{1 - EL(\Sigma_p, \widehat{\Sigma}_m^{(i)})} \right)^{-1} \\
&\stackrel{\text{a.s.}}{\rightarrow} 0,
\end{aligned}$$

if we can show further that $\frac{1}{pM} \sum_{i=1}^M \text{tr} \left((\Sigma_p - \widehat{\Sigma}_m^{(i)}) (\widehat{\Sigma}_{m,M}^{-1} - (\widehat{\Sigma}_m^{(i)})^{-1}) \right) \stackrel{\text{a.s.}}{\rightarrow} 0$, since by Theorem 5 we have $EL(\Sigma_p, \widehat{\Sigma}_m^{(i)}) \stackrel{\text{a.s.}}{\rightarrow} 0$ for each i , and that in the proof of Lemma S.4 it is proved that $p^{-1} SL(\Sigma_p, \widehat{\Sigma}_{\text{Ideal}})$ converges almost surely to a nonzero value when $\Sigma_p \neq \sigma^2 \mathbf{I}$. For this, write the term as $D_1 + D_2$, where

$$D_1 = \frac{1}{pM} \sum_{i=1}^M \text{tr}((\widehat{\Sigma}_m^{(i)} - \Sigma_p)(\widehat{\Sigma}_m^{(i)})^{-1}), \quad D_2 = \frac{1}{pM} \sum_{i=1}^M \text{tr}((\Sigma_p - \widehat{\Sigma}_m^{(i)}) \widehat{\Sigma}_{m,M}^{-1}).$$

For D_1 , writing $\mathbf{P}_{1i} = (\mathbf{p}_{1i,1}, \dots, \mathbf{p}_{1i,p})$, we can use Lemma 1 to deduce that

$$\begin{aligned}
D_1 &= \frac{1}{pM} \sum_{i=1}^M \text{tr} \left((\text{diag}(\mathbf{P}_{1i}^T \widetilde{\Sigma}_2^{(i)} \mathbf{P}_{1i}) - \mathbf{P}_{1i}^T \Sigma_p \mathbf{P}_{1i}) \text{diag}^{-1}(\mathbf{P}_{1i}^T \widetilde{\Sigma}_2^{(i)} \mathbf{P}_{1i}) \right) \\
&= \frac{1}{pM} \sum_{i=1}^M \sum_{j=1}^p \frac{\mathbf{P}_{1i,j}^T \widetilde{\Sigma}_2^{(i)} \mathbf{P}_{1i,j} - \mathbf{P}_{1i,j}^T \Sigma_p \mathbf{P}_{1i,j}}{\mathbf{P}_{1i,j}^T \widetilde{\Sigma}_2^{(i)} \mathbf{P}_{1i,j}} \\
&\leq \frac{1}{M} \sum_{i=1}^M \max_{1 \leq j \leq p} \left| \frac{\mathbf{P}_{1i,j}^T \widetilde{\Sigma}_2^{(i)} \mathbf{P}_{1i,j} - \mathbf{P}_{1i,j}^T \Sigma_p \mathbf{P}_{1i,j}}{\mathbf{P}_{1i,j}^T \widetilde{\Sigma}_2^{(i)} \mathbf{P}_{1i,j}} \right| \\
&\quad \cdot \left(1 - \max_{1 \leq j \leq p} \left| \frac{\mathbf{P}_{1i,j}^T \widetilde{\Sigma}_2^{(i)} \mathbf{P}_{1i,j} - \mathbf{P}_{1i,j}^T \Sigma_p \mathbf{P}_{1i,j}}{\mathbf{P}_{1i,j}^T \widetilde{\Sigma}_2^{(i)} \mathbf{P}_{1i,j}} \right| \right)^{-1} \stackrel{\text{a.s.}}{\rightarrow} 0.
\end{aligned}$$

For D_2 , define $\widehat{\Sigma}_{\text{CRC}}^{(i)} = \mathbf{P}_{1i} \left(\frac{1}{M} \sum_{j=1}^M \text{diag}(\mathbf{P}_{1j}^T \widetilde{\Sigma}_2^{(j)} \mathbf{P}_{1j}) \right) \mathbf{P}_{1i}$, and define also $\mathbf{B}_i = (\widehat{\Sigma}_{\text{CRC}}^{(i)} - \widehat{\Sigma}_{m,M}) (\widehat{\Sigma}_{\text{CRC}}^{(i)})^{-1}$. Then

$$\widehat{\Sigma}_{m,M}^{-1} = (\widehat{\Sigma}_{\text{CRC}}^{(i)})^{-1} (\mathbf{I} - \mathbf{B}_i)^{-1} = (\widehat{\Sigma}_{\text{CRC}}^{(i)})^{-1} + (\widehat{\Sigma}_{\text{CRC}}^{(i)})^{-1} \sum_{k \geq 1} \mathbf{B}_i^k. \quad (\text{S.4})$$

Suppose $\max_{1 \leq i, j \leq M} \|\mathbf{P}_{1i} - \mathbf{P}_{1j}\| \xrightarrow{a.s.} 0$, to be proved at the end of this proof. Then the above Neumann's series expansion is valid, since then for large enough n, p with $p/n \rightarrow c > 0$,

$$\begin{aligned} \|\mathbf{B}_i\| &\leq \frac{1}{M} \left(2 \left\| \sum_{j=1}^M (\mathbf{P}_{1i} - \mathbf{P}_{1j}) \text{diag}(\mathbf{P}_{1j}^T \widetilde{\Sigma}_2^{(j)} \mathbf{P}_{1j}) \mathbf{P}_{1i}^T \right\| \right. \\ &\quad \left. + \left\| \sum_{j=1}^M (\mathbf{P}_{1i} - \mathbf{P}_{1j}) \text{diag}(\mathbf{P}_{1j}^T \widetilde{\Sigma}_2^{(j)} \mathbf{P}_{1j}) (\mathbf{P}_{1i} - \mathbf{P}_{1j})^T \right\| \right) \\ &\quad \cdot \left(\min_{1 \leq j \leq M} \min_{1 \leq k \leq p} \mathbf{p}_{1j,k}^T \widetilde{\Sigma}_2^{(j)} \mathbf{p}_{1j,k} \right)^{-1} \\ &\leq \max_{1 \leq j \leq M} \|\mathbf{P}_{1i} - \mathbf{P}_{1j}\| \max_{1 \leq k \leq p} \mathbf{p}_{1j,k}^T \Sigma_p \mathbf{p}_{1j,k} \left(2 + \|\mathbf{P}_{1i} - \mathbf{P}_{1j}\| \right) \\ &\quad \cdot \left(\min_{1 \leq j \leq M} \min_{1 \leq k \leq p} \mathbf{p}_{1j,k}^T \Sigma_p \mathbf{p}_{1j,k} \right)^{-1} \\ &\leq \max_{1 \leq i, j \leq M} \|\mathbf{P}_{1i} - \mathbf{P}_{1j}\| \|\Sigma_p\| (2 + \|\mathbf{P}_{1i} - \mathbf{P}_{1j}\|) / \lambda_{\min}(\Sigma_p) \xrightarrow{a.s.} 0, \quad (\text{S.5}) \end{aligned}$$

where the second last line follows from Lemma 1 when n, p is large. With (S.4), $D_2 = D_3 + D_4$, where

$$\begin{aligned} D_3 &= \frac{1}{pM} \sum_{i=1}^M \text{tr} \left((\Sigma_p - \widehat{\Sigma}_m^{(i)}) (\widehat{\Sigma}_{\text{CRC}}^{(i)})^{-1} \right), \\ D_4 &= \frac{1}{pM} \sum_{i=1}^M \text{tr} \left((\Sigma_p - \widehat{\Sigma}_m^{(i)}) (\widehat{\Sigma}_{\text{CRC}}^{(i)})^{-1} \sum_{k \geq 1} \mathbf{B}_i^k \right). \end{aligned}$$

Similar to the analysis of D_1 , using Lemma 1, we have

$$\begin{aligned} D_3 &= \frac{1}{pM} \sum_{i=1}^M \sum_{k=1}^p \frac{\mathbf{p}_{1i,k}^T \Sigma_p \mathbf{p}_{1i,k} - \mathbf{p}_{1i,k}^T \widetilde{\Sigma}_2^{(i)} \mathbf{p}_{1i,k}}{\frac{1}{M} \sum_{j=1}^M \mathbf{p}_{1j,k}^T \widetilde{\Sigma}_2^{(j)} \mathbf{p}_{1j,k}} \\ &\leq \max_{1 \leq i \leq M} \max_{1 \leq k \leq p} \left| \frac{\mathbf{p}_{1i,k}^T \Sigma_p \mathbf{p}_{1i,k} - \mathbf{p}_{1i,k}^T \widetilde{\Sigma}_2^{(i)} \mathbf{p}_{1i,k}}{\mathbf{p}_{1i,k}^T \Sigma_p \mathbf{p}_{1i,k}} \right| \cdot \mathbf{p}_{1i,k}^T \Sigma_p \mathbf{p}_{1i,k} \\ &\quad \cdot \left(\min_{1 \leq j \leq M} \min_{1 \leq k \leq p} \left(\frac{\mathbf{p}_{1j,k}^T \widetilde{\Sigma}_2^{(j)} \mathbf{p}_{1j,k} - \mathbf{p}_{1j,k}^T \Sigma_p \mathbf{p}_{1j,k}}{\mathbf{p}_{1j,k}^T \Sigma_p \mathbf{p}_{1j,k}} \right) \right) \mathbf{p}_{1j,k}^T \Sigma_p \mathbf{p}_{1j,k} \\ &\quad \left. + \mathbf{p}_{1j,k}^T \Sigma_p \mathbf{p}_{1j,k} \right)^{-1} \xrightarrow{a.s.} 0. \end{aligned}$$

For D_4 , with $\|\mathbf{B}_i\| \xrightarrow{a.s.} 0$ proved in (S.5) and using Lemma 1, we have

$$\begin{aligned} D_4 &\leq \max_{1 \leq i \leq M} \|\widehat{\boldsymbol{\Sigma}}_m^{(i)} - \boldsymbol{\Sigma}_p\| \cdot \lambda_{\min}^{-1}(\widehat{\boldsymbol{\Sigma}}_{\text{CRC}}^{(i)}) \cdot \sum_{k \geq 1} \|\mathbf{B}_i\|^k \\ &\leq \frac{(\|\boldsymbol{\Sigma}_p\| + \max_{1 \leq i \leq M} \max_{1 \leq k \leq p} \mathbf{p}_{1i,k}^T \widetilde{\boldsymbol{\Sigma}}_2^{(i)} \mathbf{p}_{1i,k}) \|\mathbf{B}_i\|}{\min_{1 \leq i \leq M} \min_{1 \leq k \leq p} \mathbf{p}_{1i,k}^T \widetilde{\boldsymbol{\Sigma}}_2^{(i)} \mathbf{p}_{1i,k} (1 - \|\mathbf{B}_i\|)} \xrightarrow{a.s.} 0. \end{aligned}$$

Hence if we can show further that $\max_{1 \leq i, j \leq M} \|\mathbf{P}_{1i} - \mathbf{P}_{1j}\| \xrightarrow{a.s.} 0$, the proof of the Theorem is complete. To this end, using (S.3), we have for any $1 \leq i \leq M$, with $\lambda_{1i,k}$ denoting the k th largest eigenvalue of $\widetilde{\boldsymbol{\Sigma}}_1^{(i)}$,

$$p^{-1} \sum_{k=1}^p g(\mathbf{p}_{1i,k}^T \boldsymbol{\Sigma}_p \mathbf{p}_{1i,k}) \mathbf{1}_{\{\lambda_{1i,k} \leq x\}} \xrightarrow{a.s.} \int_{-\infty}^x g(\delta(\lambda)) dF(\lambda),$$

since $m/n \rightarrow 1$ implies that $p/n, p/m \rightarrow c > 0$. Setting $g \equiv 1$, it is easy to see that $\lambda_{1i,k}$ is almost surely the same as $\lambda_{1j,k}$ for $1 \leq i, j \leq M$ and $1 \leq k \leq p$, as $n, p \rightarrow \infty$. These imply that for $1 \leq i \leq M$, $\mathbf{p}_{1i,k}^T \boldsymbol{\Sigma}_p \mathbf{p}_{1i,k}$ is almost surely the same as $\delta(\lambda_{11,k})$, and hence for $1 \leq i, j \leq M$ and $1 \leq k \leq p$,

$$\mathbf{p}_{1i,k}^T \boldsymbol{\Sigma}_p \mathbf{p}_{1i,k} - \mathbf{p}_{1j,k}^T \boldsymbol{\Sigma}_p \mathbf{p}_{1j,k} \xrightarrow{a.s.} 0.$$

But $\mathbf{p}_{1i,k}^T \boldsymbol{\Sigma}_p \mathbf{p}_{1i,k} - \mathbf{p}_{1j,k}^T \boldsymbol{\Sigma}_p \mathbf{p}_{1j,k} = (\mathbf{p}_{1i,k} - \mathbf{p}_{1j,k})^T \boldsymbol{\Sigma}_p (\mathbf{p}_{1i,k} + \mathbf{p}_{1j,k})$, hence if $\mathbf{p}_{1i,k}$ or $\mathbf{p}_{1j,k}$ are not eigenvectors of $\boldsymbol{\Sigma}_p$ asymptotically (this excludes the case $\boldsymbol{\Sigma}_p = \sigma^2 \mathbf{I}$, which by assumption is not the case), then there exists a constant $a > 0$ such that

$$(\mathbf{p}_{1i,k} - \mathbf{p}_{1j,k})^T \boldsymbol{\Sigma}_p (\mathbf{p}_{1i,k} + \mathbf{p}_{1j,k}) = \|\mathbf{p}_{1i,k} - \mathbf{p}_{1j,k}\| \cdot a \|\mathbf{p}_{1i,k} + \mathbf{p}_{1j,k}\| \xrightarrow{a.s.} 0,$$

implying that for all $1 \leq i, j \leq M$, $1 \leq k \leq p$, $\|\mathbf{p}_{1i,k} \pm \mathbf{p}_{1j,k}\| \xrightarrow{a.s.} 0$. With the correct orientation chosen, we must then have $\max_{1 \leq i, j \leq M} \|\mathbf{P}_i - \mathbf{P}_j\| \xrightarrow{a.s.} 0$.

Finally, if both $\mathbf{p}_{1i,k}$ and $\mathbf{p}_{1j,k}$ are asymptotic eigenvectors of $\boldsymbol{\Sigma}_p$ for some $1 \leq k \leq p$, then assuming $\|\boldsymbol{\Sigma}_p \mathbf{p}_{1i,k} - \lambda \mathbf{p}_{1i,k}\| \xrightarrow{a.s.} 0$ and $\|\boldsymbol{\Sigma}_p \mathbf{p}_{1j,k} - \gamma \mathbf{p}_{1j,k}\| \xrightarrow{a.s.} 0$,

$$\begin{aligned} (\mathbf{p}_{1i,k} - \mathbf{p}_{1j,k})^T \boldsymbol{\Sigma}_p (\mathbf{p}_{1i,k} + \mathbf{p}_{1j,k}) &= (\mathbf{p}_{1i,k} - \mathbf{p}_{1j,k})^T (\lambda \mathbf{p}_{1i,k} + \gamma \mathbf{p}_{1j,k}) \\ &\quad + (\mathbf{p}_{1i,k} - \mathbf{p}_{1j,k})^T (\boldsymbol{\Sigma}_p \mathbf{p}_{1i,k} - \lambda \mathbf{p}_{1i,k}) \\ &\quad + \boldsymbol{\Sigma}_p \mathbf{p}_{1j,k} - \gamma \mathbf{p}_{1j,k}) \\ &= (\lambda - \gamma)(1 - \mathbf{p}_{1i,k}^T \mathbf{p}_{1j,k}) + o(1), \end{aligned}$$

so that if $\lambda \neq \gamma$, we must have $\|\mathbf{p}_{1i,k} - \mathbf{p}_{1j,k}\| \xrightarrow{a.s.} 0$, which is what we want. If $\lambda = \gamma$, then $\mathbf{p}_{1i,k}$ and $\mathbf{p}_{1j,k}$ are asymptotically spanning the same eigenspace, hence $\max_{1 \leq i, j \leq M} \|\mathbf{P}_i - \mathbf{P}_j\| \xrightarrow{a.s.} 0$ still holds. This completes the proof of the theorem. \square

References

- Abadir, K. M., Distaso, W., and Žikeš, F. (2014). Design-free estimation of variance matrices. *Journal of Econometrics*, 181(2):165 – 180.
- Bai, Z. and Silverstein, J. (2010). *Spectral Analysis of Large Dimensional Random Matrices*. Springer Series in Statistics, New York, 2 edition.
- Bai, Z. D. and Silverstein, J. (1998). No eigenvalues outside the support of the limiting spectral distribution of large-dimensional sample covariance matrices. *The Annals of Probability*, 26(1):316–345.
- Fan, J., Liao, Y., and Mincheva, M. (2013). Large covariance estimation by thresholding principal orthogonal complements. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 75(4):603–680.
- Friedman, J., Hastie, T., and Tibshirani, R. (2008). Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441.
- Ledoit, O. and Péché, S. (2011). Eigenvectors of some large sample covariance matrix ensembles. *Probability Theory and Related Fields*, 151(1-2):233–264.
- Ledoit, O. and Wolf, M. (2012). Nonlinear shrinkage estimation of large-dimensional covariance matrices. *The Annals of Statistics*, 40(2):1024–1060.
- Ledoit, O. and Wolf, M. (2013a). Optimal estimation of a large-dimensional covariance matrix under Stein’s loss. ECON - Working Papers 122, Department of Economics - University of Zurich.
- Ledoit, O. and Wolf, M. (2013b). Spectrum estimation: a unified framework for covariance matrix estimation and PCA in large dimensions. ECON - Working Papers 105, Department of Economics - University of Zurich.
- Silverstein, J. and Bai, Z. (1995). On the empirical distribution of eigenvalues of a class of large dimensional random matrices. *Journal of Multivariate Analysis*, 54(2):175 – 192.