



Published in final edited form as:

Stat Sin. 2021 April ; 31(2): 571–601. doi:10.5705/ss.202019.0052.

Asymptotics of eigenstructure of sample correlation matrices for high-dimensional spiked models

David Morales-Jimenez¹, Iain M. Johnstone², Matthew R. McKay³, Jeha Yang²

¹ECIT Institute, Queen's University Belfast, UK

²Department of Statistics, Stanford University, USA

³ECE Department, Hong Kong University of Science and Technology, Hong Kong

Abstract

Sample correlation matrices are widely used, but for high-dimensional data little is known about their spectral properties beyond “null models”, which assume the data have independent coordinates. In the class of spiked models, we apply random matrix theory to derive asymptotic first-order and distributional results for both leading eigenvalues and eigenvectors of sample correlation matrices, assuming a high-dimensional regime in which the ratio p/n , of number of variables p to sample size n , converges to a positive constant. While the first-order spectral properties of sample correlation matrices match those of sample covariance matrices, their asymptotic distributions can differ significantly. Indeed, the correlation-based fluctuations of both sample eigenvalues and eigenvectors are often remarkably smaller than those of their sample covariance counterparts.

Keywords

Sample correlation; eigenstructure; spiked models

1. Introduction

Estimating a correlation matrix is a fundamental statistical task. It is widely applied in areas such as viral sequence analysis and vaccine design in biology (Dahirel et al., 2011, Quadeer et al., 2014, 2018), large portfolio design in finance (Plerou et al., 2002), signal detection in radio astronomy (Leshem and van der Veen, 2001), and collaborative filtering (Liu et al., 2014, Ruan et al., 2016), among many others. In classical statistical settings, with a limited number of variables p and a large sample size n , the sample correlation matrix performs well and its statistical properties are well understood; see, for example, Girshick (1939), Konishi (1979), Fang and Krishnaiah (1982), Schott (1991), Kollo and Neudecker (1993), and Boik (2003). Modern applications, however, often exhibit high dimensionality, with large p and,

David Morales-Jimenez: ECIT Institute, Queen's University Belfast, UK, d.morales@qub.ac.uk.

Supplementary Material

The online Supplementary Material provides proofs for the following: (i) the Gaussian particularizations of our main results (Corollaries 1 and 2); (ii) the instrumental tightness properties in Lemma 3; and (iii) the asymptotic properties of normalized bilinear forms in Lemma 1 and Proposition 1; see Sections S1, S2, and S3, respectively.

in many cases, limited n . In such cases, sample correlation matrices become inaccurate owing to an aggregation of statistical noise across the matrix coordinates that is visible in the eigen-spectrum (El Karoui, 2009). This is particularly important in principal component analysis (PCA), which often involves projecting data onto the leading eigenvectors of the sample correlation matrix or, equivalently, onto those of the sample covariance matrix after standardizing the data.

Despite the extensive use of sample correlation matrices, relatively little is known about theoretical properties of their eigen-spectra in high dimensions. In contrast, sample covariance matrices have been studied extensively, and a rich body of literature now exists (e.g., Yao et al. (2015)). Their asymptotic properties have typically been described in high-dimensional settings in which the number of samples and variables both grow large, often though not always at the same rate, based on the theory of random matrices. Specific first- and second-order results for the eigenvalues and eigenvectors of sample covariance matrices are reviewed in Bai and Silverstein (2009), Couillet and Debbah (2011), and Yao et al. (2015).

For the spectra of high-dimensional sample *correlation* matrices, current theoretical results focus on the simplest “null model” scenario, in which the data are assumed to be independent. In this null model, correlation matrices share many of the same asymptotic properties as covariance matrices from independent and identically distributed (i.i.d.) data, with zero mean and unit variance. Thus, the empirical eigenvalue distribution converges to the Marchenko–Pastur distribution, almost surely (Jiang, 2004b), and the largest and smallest eigenvalues converge to the edges of this distribution (Jiang, 2004b, Xiao and Zhou, 2010). Moreover, the rescaled largest and smallest eigenvalues asymptotically follow the Tracy–Widom law (Bao et al., 2012, Pillai and Yin, 2012). Central limit theorems (CLTs) for linear spectral statistics have also been derived (Gao et al., 2017). A separate line of work studies the maximum absolute off-diagonal entry of sample correlation matrices, referred to as “coherence” (Jiang, 2004a, Cai and Jiang, 2011, 2012), which has been proposed as a statistic for conducting independence tests; see also Cochran et al. (1995), Mestre and Vallet (2017), and the references therein. Hero and Rajaratnam (2011, 2012) use a related statistic to identify variables exhibiting strong correlations, an approach referred to as “correlation screening.”

For non-trivial correlation models, however, asymptotic results for the spectra of sample correlation matrices are quite scarce. Notably, El Karoui (2009) shows that, for a fairly general class of covariance models with bounded spectral norm, to first order, the eigenvalues of sample correlation matrices asymptotically coincide with those of sample covariance matrices with unit-variance data, generalizing earlier results of Jiang (2004b) and Xiao and Zhou (2010). Under similar covariance assumptions, recent work also presents CLTs for linear spectral statistics of sample correlation matrices (Mestre and Vallet, 2017), extending the work of Gao et al. (2017). First order behavior again coincides with that of sample covariances. However, the asymptotic fluctuations are quite different for sample correlation matrices.

This study considers a particular class of correlation matrix models, the so-called “spiked models,” in which a few large or small eigenvalues of the population covariance (or correlation) matrix are assumed to be well separated from the rest (Johnstone, 2001). Spiked covariance models are relevant in applications in which the primary covariance information lies in a relatively small number of eigenmodes. Such applications include collaborative signal detection in cognitive radio systems (Bianchi et al., 2009), fault detection in sensor networks (Couillet and Hachem, 2013), adaptive beamforming in array processing (Hachem et al., 2013, Vallet et al., 2015, Yang et al., 2018), and protein contact prediction in biology (Cocco et al., 2011, 2013). The spectral properties of spiked covariance models have been well studied, with precise analytical results established for the asymptotic first-order and distributional properties of both eigenvalues and eigenvectors; see, for example, Baik et al. (2005), Baik and Silverstein (2006), Paul (2007), Bai and Yao (2008), Benaych-Georges and Nadakuditi (2011), Couillet and Hachem (2013), Bloemendal et al. (2016). For reviews, see also Couillet and Debbah (2011, Chapter 9) and Yao et al. (2015, Chapter 11).

Less is known about the spectrum of sample correlation matrices under spiked models. Although the asymptotic first-order behavior is expected to coincide with that of the sample covariance, as a consequence of El Karoui (2009), a simple simulation reveals striking differences in the fluctuations of both sample eigenvalues and eigenvectors; see Figure 1.

Here, we present theoretical results to describe these observed phenomena. We obtain asymptotic first-order and distribution results for the eigenvalues and eigenvectors of sample correlation matrices under a spiked model. Paul (2007) proved theorems for sample covariance matrices in the special case of Gaussian data. In essence, we present analogs of these theorems for sample correlation matrices, and extend them to non-Gaussian data. To first order, the eigenvalues and eigenvectors coincide asymptotically with those of sample covariance matrices; however, their fluctuations can be very different. Indeed, for both the largest sample correlation eigenvalues (Theorem 1) and the projections of the corresponding eigenvectors (Theorem 2), the asymptotic variances admit a decomposition into three terms. The first term is just the asymptotic variance for sample covariance matrices generated from Gaussian data; the second adds corrections due to non-Gaussianity, and the third captures further corrections due to data normalization imposed by the sample correlation matrix. (This last amounts to normalizing the entries of the sample covariance matrix using the sample variances). Consistent with the example shown in Figure 1(a), in the CLT for the leading sample eigenvalues, the sample correlation eigenvalues often show lower fluctuations—despite the variance normalization—than those of the sample covariance eigenvalues. As seen in Figure 1(b), the (normalized) eigenvector projections are typically asymptotically correlated, even for Gaussian data, unlike the sample covariance setting of Paul (2007, Theorem 5).

Technical contributions

We build on and extend a set of random matrix tools for studying spiked covariance models. The companion manuscript (Johnstone and Yang, 2018) [JY], gives an exposition and parallel treatment for sample covariance matrices. Important adaptations are needed here to account for the data normalization imposed by sample correlation matrices. Among key

technical contributions of our work, basic to our main theorems, are asymptotic first-order and distributional properties for bilinear forms and matrix quadratic forms with normalized entries, Section 4. A novel regularization-based proof strategy is used to establish the inconsistency of eigenvector projections in the case of “subcritical” spiked eigenvalues, Theorem 3.

Model M

Let $x \in \mathbb{R}^{m+p}$ be a random vector with finite $(4+\delta)$ th moment for some $\delta > 0$. Consider the partition

$$x = \begin{bmatrix} \xi \\ \eta \end{bmatrix}.$$

Assume that $\xi \in \mathbb{R}^m$ has mean zero and covariance Σ , and is independent of $\eta \in \mathbb{R}^p$, which has i.i.d components η_i with mean zero and unit variance. Let $\Sigma_D = \text{diag}(\sigma_1^2, \dots, \sigma_m^2)$ be the diagonal matrix containing the variances of ξ_i , and let $\Gamma = \Sigma_D^{-1/2} \Sigma \Sigma_D^{-1/2}$ be the correlation matrix of ξ with eigen-decomposition $\Gamma = PLP^T$, where $P = [p_1, \dots, p_m]$ is the eigenvector matrix, and $L = \text{diag}(\ell_1, \dots, \ell_m)$ contains the spike correlation eigenvalues $\ell_1 \dots \ell_m > 0$.

The correlation matrix of x is therefore $\Gamma_x = \text{blkdiag}(\Gamma, I)$, with eigenvalues $\ell_1, \dots, \ell_m, 1, \dots, 1$, and corresponding eigenvectors $\mathbf{p}_1, \dots, \mathbf{p}_m, e_{m+1}, \dots, e_{m+p}$ where $\mathbf{p}_i = [p_i^T \ 0_p^T]^T$ and e_j is the j th canonical vector (i.e., a vector of all zeros, except for a one in the j th coordinate).

Consider a sequence of i.i.d. copies of x , the first n of which fill the columns of the $(m+p) \times n$ data matrix $X = (x_{ij})$. We assume m is fixed, whereas p and n increase with

$$\gamma_n = p/n \rightarrow \gamma > 0 \quad \text{as } p, n \rightarrow \infty.$$

Notation

Let $S = n^{-1}XX^T$ be the sample covariance matrix, and $S_D = \text{diag}(\hat{\sigma}_1^2, \dots, \hat{\sigma}_{m+p}^2)$ be the diagonal matrix containing the sample variances. Let $R = S_D^{-1/2} S S_D^{-1/2}$ be the sample correlation matrix, with corresponding ν th sample eigenvalue and eigenvector satisfying

$$R\hat{\mathbf{p}}_\nu = \hat{\lambda}_\nu \hat{\mathbf{p}}_\nu,$$

where, for later use, we partition $\hat{\mathbf{p}}_\nu = [\hat{p}_\nu^T, \hat{v}_\nu^T]^T$. Here \hat{p}_ν is the subvector of $\hat{\mathbf{p}}_\nu$ restricted to the first m coordinates.

For $\ell > 1 + \sqrt{\gamma}$, define

$$\rho(\ell, \gamma) = \ell + \gamma \frac{\ell}{\ell - 1}, \quad \dot{\rho}(\ell, \gamma) = \frac{\partial \rho(\ell, \gamma)}{\partial \ell} = 1 - \frac{\gamma}{(\ell - 1)^2}.$$

For an index ν , for which $\ell_\nu > 1 + \sqrt{\gamma}$ is a simple eigenvalue, set

$$\rho_\nu = \rho(\ell_\nu, \gamma), \quad \rho_{\nu n} = \rho(\ell_\nu, \gamma_n), \quad \dot{\rho}_\nu = \dot{\rho}(\ell_\nu, \gamma), \quad \dot{\rho}_{\nu n} = \dot{\rho}(\ell_\nu, \gamma_n). \quad (1.1)$$

We refer to eigenvalues satisfying $\ell_\nu > 1 + \sqrt{\gamma}$ as “supercritical,” and those satisfying $\ell_\nu \leq 1 + \sqrt{\gamma}$ as “subcritical,” with the quantity $1 + \sqrt{\gamma}$ referred to as the “phase transition.”

To describe and interpret the variance terms in the limiting distributions to follow, we need some definitions. Let $\bar{\xi}_i = \xi_i/\sigma_i$ and $\kappa_{ij} = \mathbb{E}\bar{\xi}_i\bar{\xi}_j$ denote the scaled components of ξ and their covariances; of course $\kappa_{ii} = 1$. The corresponding scaled fourth-order cumulants are

$$\kappa_{ijj'j'} = \mathbb{E}[\bar{\xi}_i\bar{\xi}_j\bar{\xi}_{i'}\bar{\xi}_{j'}] - \kappa_{ij}\kappa_{i'j'} - \kappa_{ij'}\kappa_{ji'} - \kappa_{ii'}\kappa_{jj'}. \quad (1.2)$$

When ξ is Gaussian, $\kappa_{ijj'j'} \equiv 0$.

The effect of variance scaling in the correlation matrix is described using additional quadratic functions of $(\bar{\xi}_i)$, defined by

$$\chi_{ij} = \bar{\xi}_i\bar{\xi}_j, \quad \psi_{ij} = \kappa_{ij}(\bar{\xi}_i^2 + \bar{\xi}_j^2)/2 \quad (1.3)$$

$$\check{\kappa}_{ijj'j'} = \text{Cov}(\psi_{ij}, \psi_{i'j'}) - \text{Cov}(\psi_{ij}, \chi_{i'j'}) - \text{Cov}(\chi_{ij}, \psi_{i'j'}). \quad (1.4)$$

Tensor notation

For convenience, it is useful to consider $\kappa_{ijj'j'}$ and $\check{\kappa}_{ijj'j'}$ as entries of four-dimensional tensor arrays κ and $\check{\kappa}$, respectively, and to define an additional array $\mathcal{P}^{\mu\mu'vv'}$ with entries $p_{\mu, i} p_{\mu', j} p_{\nu, i'} p_{\nu', j'}$. In addition, define \mathcal{P}^ν as $\mathcal{P}^{\nu\nu\nu\nu}$. Finally, for a second array A of the same dimensions,

$$[\mathcal{P}^\nu, A] = \sum_{i, j, i', j'} P_{ijj'j'}^\nu A_{ijj'j'}.$$

2. Main results

Our first main result, proved in Section 5, gives the asymptotic properties of the largest (spike) eigenvalues of the sample correlation matrix:

Theorem 1

Assume Model M, and that $\ell_\nu > 1 + \sqrt{\gamma}$ is a simple eigenvalue. As $p/n \rightarrow \gamma > 0$,

$$\begin{aligned} (i) \quad & \widehat{\ell}_\nu \xrightarrow{\text{a.s.}} \rho_\nu, \\ (ii) \quad & \sqrt{n}(\widehat{\ell}_\nu - \rho_{\nu n}) \xrightarrow{\mathcal{D}} N(0, \bar{\sigma}_\nu^2), \end{aligned} \quad (2.5)$$

where

$$\tilde{\sigma}_v^2 = 2\dot{\rho}_v \ell_v^2 + \dot{\rho}_v^2[\mathcal{P}^v, \kappa] + \dot{\rho}_v^2[\mathcal{P}^v, \tilde{\kappa}]. \tag{2.6}$$

Centering at ρ_{vn} rather than at ρ_v is important. If, for example, $\gamma_n = \gamma + an^{-1/2}$, then

$$\sqrt{n}(\hat{\ell}_v - \rho_v) \xrightarrow{\mathcal{D}} N(a\ell_v(\ell_v - 1)^{-1}, \tilde{\sigma}_v^2),$$

and we see a limiting shift. Furthermore, it may also be beneficial to consider $\tilde{\sigma}_{vn}^2$ instead of $\tilde{\sigma}_v^2$, obtained by replacing $\dot{\rho}_v$ with $\dot{\rho}_{vn}$ in (2.6), such that

$$\sqrt{n}(\hat{\ell}_v - \rho_{vn})/\tilde{\sigma}_{vn} \xrightarrow{\mathcal{D}} N(0, 1).$$

The asymptotic first-order limit in (j), which follows as an easy consequence of El Karoui (2009), coincides with that of the v th largest eigenvalue of a sample covariance matrix computed from data with population covariance Γ (Paul, 2007). This implies that, when constructing \mathcal{R} , normalizing by the sample variances has no effect on the leading eigenvalues, at least to first order.

However, key differences are seen when looking at the asymptotic distribution, given in (ii), and in the variance formula (2.6) in particular. This can be readily interpreted. The first term corresponds to the variance in the Gaussian-covariance case of Paul (2007), again for samples with covariance Γ . The second provides a correction of that result for non-Gaussian data, see the companion article [JY]. The third term describes the contribution specific to sample correlation matrices, representing the effect of normalizing the data by the sample variances. This term is often negative, and is evaluated explicitly for Gaussian data in Corollary 1 below, proved in the Supplementary Material, S1.1.

Corollary 1

For ξ Gaussian, the asymptotic variance in Theorem 1 simplifies to

$$\tilde{\sigma}_v^2 = 2\ell_v^2 \dot{\rho}_v \left[1 - \dot{\rho}_v \left(2\ell_v \text{tr} P_{D,v}^4 - \text{tr}(P_{D,v} \Gamma P_{D,v}) \right)^2 \right],$$

where $P_{D,v} = \text{diag}(p_{v,1}, \dots, p_{v,m})$.

Thus, computing the sample correlation results in the asymptotic variance being scaled by $1 - \dot{\rho}_v \Delta_v$, relative to the sample covariance, where

$$\Delta_v = 2\ell_v \text{tr} P_{D,v}^4 - \text{tr}(P_{D,v} \Gamma P_{D,v})^2 = 2\ell_v \sum_i p_{v,i}^4 - \sum_{i,j} (p_{v,i} \kappa_{ij} p_{v,j})^2$$

is often positive, implying that spiked eigenvalues of the sample correlation often exhibit a smaller variance than those of the sample covariance. Indeed, such variance reduction occurs iff

$$\sum_{i,j} (p_{\nu, i} \kappa_{ij} p_{\nu, j})^2 < 2\ell_{\nu} \sum_i p_{\nu, i}^4 = \sum_{i,j} p_{\nu, i} \kappa_{ij} p_{\nu, j} (p_{\nu, i}^2 + p_{\nu, j}^2), \tag{2.7}$$

with the last identity following from the fact that $\ell_{\nu} p_{\nu, i} = \sum_j \kappa_{ij} p_{\nu, j}$. Condition (2.7), and variance reduction, holds in the following cases:

- i. both Γ and p_{ν} have nonnegative entries, or
- ii. $2\ell_{\nu} \sum_i p_{\nu, i}^4 > 1$, or
- iii. $2\ell_{\nu} > \ell_1^2$.

In case (i), the inequalities $0 \leq p_{\nu, i} \kappa_{ij} p_{\nu, j} \leq 2p_{\nu, i} p_{\nu, j} \leq p_{\nu, i}^2 + p_{\nu, j}^2$ yield (2.7). Note that if Γ has nonnegative entries, then the Perron–Frobenius theorem establishes the existence of an eigenvector with nonnegative components for ℓ_1 ; furthermore, if Γ has positive entries, by the same theorem, ℓ_1 is simple and associated with an eigenvector with positive components.

Case (ii) follows from $\sum_{i,j} (p_{\nu, i} \kappa_{ij} p_{\nu, j})^2 \leq \sum_{i,j} (p_{\nu, i} p_{\nu, j})^2 = 1$, and holds if $\ell_{\nu} > m/2$, because $\sum_i p_{\nu, i}^4 \geq 1/m$. Case (iii) follows from the inequalities $2p_{\nu, i}^2 p_{\nu, j}^2 \leq p_{\nu, i}^4 + p_{\nu, j}^4$ and $\sum_j \kappa_{ij}^2 = (\Gamma^2)_{ii} \leq \|\Gamma^2\| = \ell_1^2$. Note that this is rather special, in that it has nothing to do with eigenvectors, and a necessary condition for it to hold is $\ell_1 > 2$.

Condition (2.7) can fail, however. For example, for even m and $r \in (0, 1)$, consider

$$\Gamma = \begin{pmatrix} 1 & -r \\ -r & 1 \end{pmatrix} \otimes 1_{m/2} 1_{m/2}^T,$$

where $1_{m/2}$ is the $(m/2)$ -dimensional vector of all ones, which corresponds to two negatively correlated groups of identical random vectors. This has simple supercritical eigenvalues $\ell_1 = (1+r)m/2$ and $\ell_2 = (1-r)m/2$ when $m > 2(1+\sqrt{\gamma})(1-r)$, with $p_{\nu, i}^2 = m^{-1}$ for $\nu = 1, 2$. One finds that $\ell_2 = (1-2r-r^2)/2 < 0$ for $r > \sqrt{2}-1$, although $\ell_1 > 0$ because $\ell_1 > m/2$, which implies case (ii).

We turn now to the eigenvectors. Again, fix an index ν for which $\ell_{\nu} > 1 + \sqrt{\gamma}$ is a simple eigenvalue of Γ , with corresponding eigenvector $\mathbf{p}_{\nu} = [p_{\nu}^T \ 0_p^T]^T$. Recall that $\hat{\mathbf{p}}_{\nu} = [\hat{p}_{\nu}^T \ \hat{v}_{\nu}^T]^T$ is the ν th sample eigenvector of R , and let $a_{\nu} = \hat{p}_{\nu} / \|\hat{p}_{\nu}\|$ be the corresponding normalized subvector of $\hat{\mathbf{p}}_{\nu}$, restricted to the first m coordinates. The next result establishes a limit for the eigenvector projection $\langle \hat{\mathbf{p}}_{\nu}, \mathbf{p}_{\nu} \rangle$, and a CLT for the normalized cross-projections $P^T a_{\nu} = [p_1^T a_{\nu}, \dots, p_m^T a_{\nu}]^T$; see Sections 6.1 and 6.2.

Theorem 2

Assume Model M, and that $\ell_{\nu} > 1 + \sqrt{\gamma}$ is a simple eigenvalue. Then, as $p/n \rightarrow \gamma > 0$,

$$(i) \quad \langle \hat{\mathbf{p}}_{\nu}, \mathbf{p}_{\nu} \rangle^2 \xrightarrow{\text{a.s.}} \hat{\rho}_{\nu} \ell_{\nu} / \rho_{\nu},$$

$$(ii) \quad \sqrt{n}(\mathbf{P}^T a_{\nu} - e_{\nu}) \xrightarrow{\mathcal{D}} N(0, \Sigma_{\nu}),$$

where $\Sigma_{\nu} = \mathcal{D}_{\nu} \tilde{\Sigma}_{\nu} \mathcal{D}_{\nu}$ with

$$\mathcal{D}_{\nu} = \sum_{k \neq \nu}^m (\ell_{\nu} - \ell_k)^{-1} e_k e_k^T \quad (2.8)$$

$$\tilde{\Sigma}_{\nu, kl} = \hat{\rho}_{\nu}^{-1} \ell_k \ell_{\nu} \delta_{k,l} + [\mathcal{P}^{k\nu l\nu}, \kappa] + [\mathcal{P}^{k\nu l\nu}, \check{\kappa}], \quad (2.9)$$

where $\delta_{k,l} = 1$ if $k = l$, and zero otherwise.

The CLT result in (ii) can be rephrased in terms of the entries of a_{ν} , for which we readily obtain $\sqrt{n}(a_{\nu} - p_{\nu}) \xrightarrow{\mathcal{D}} N(0, P\Sigma_{\nu}P^T)$; note that Σ_{ν} has zeros in the ν th row and the ν th column.

As for the eigenvalues, Theorem 2 shows that the spiked eigenvectors of sample correlation matrices exhibit the same first-order behavior as those of the sample covariance (Paul, 2007). The difference again lies in the asymptotic fluctuations, captured by the covariance matrix Σ_{ν} . Note that this is decomposed as a product of \mathcal{D}_{ν} —a diagonal matrix—and the matrix $\tilde{\Sigma}_{\nu}$, which involves the three terms in (2.9). These terms have similar interpretations as those discussed previously in (2.6). That is, the first term captures the asymptotic fluctuations for a Gaussian-covariance model (Paul, 2007), the second term captures the effect of non-Gaussianity in the covariance case [JY], and the third term captures information specific to the correlation case, representing fluctuations due to sample variance normalization. Note that only the first term is diagonal in general, suggesting that the eigenvector projections may be asymptotically correlated, as seen earlier in Figure 1(b), right panel. This holds also for Gaussian data, evaluated explicitly in Corollary 2 below; see Supplementary Material, S1.2, for the proof. We note an interesting contrast with the eigenvector projections for covariance matrices (Paul, 2007), described only by the leading term in (2.9).

Corollary 2

For ξ Gaussian, the asymptotic covariance in Theorem 2 reduces to $\Sigma_{\nu} = \mathcal{D}_{\nu} \tilde{\Sigma}_{\nu} \mathcal{D}_{\nu}$,

$$\tilde{\Sigma}_{\nu} = \frac{\ell_{\nu}}{\rho_{\nu}} L + (\ell_{\nu} I + L) \left(\frac{1}{2} \mathcal{X} - \ell_{\nu} \mathcal{Y} \right) (\ell_{\nu} I + L) + \ell_{\nu} \left(\ell_{\nu}^2 \mathcal{Y} - L \mathcal{Y} L \right),$$

where $\mathcal{X} = P^T P_{D, \nu} (\Gamma \circ \Gamma) P_{D, \nu} P$, $\mathcal{Y} = P^T P_{D, \nu}^2 P$, and \circ denotes the Hadamard product.

Thus, for Gaussian data, the entries of the asymptotic covariance matrix are given by (for $k, l \neq \nu$)

$$\Sigma_{v,kl} = (\ell_v - \ell_k)^{-1} (\ell_v - \ell_l)^{-1} \left[\frac{\ell_v}{\rho_v} \ell_k \delta_{k,l} + (\ell_v + \ell_k)(\ell_v + \ell_l) \frac{\mathcal{F}_{kl}}{2} - \ell_v(\ell_v(\ell_k + \ell_l) + 2\ell_k \ell_l) \mathcal{Y}_{kl} \right].$$

Consider now the subcritical case in which v is such that $1 < \ell_v \leq 1 + \sqrt{\gamma}$. Let \mathbf{p}_v denote the corresponding population eigenvector, and let $\hat{\ell}_v$ and $\hat{\mathbf{p}}_v$ denote the corresponding sample eigenvalue and eigenvector, respectively. With proofs deferred to Sections 5.1 and 6.3, we have the following result:

Theorem 3

Assume Model M, and that $1 < \ell_v \leq 1 + \sqrt{\gamma}$ is a simple eigenvalue. Then, as $p/n \rightarrow \gamma > 0$,

- (i) $\hat{\ell}_v \xrightarrow{\text{a.s.}} (1 + \sqrt{\gamma})^2,$
- (ii) $\langle \hat{\mathbf{p}}_v, \mathbf{p}_v \rangle^2 \xrightarrow{\text{a.s.}} 0.$

Once again, the asymptotic first-order limits of the sample eigenvalue and its associated eigenvector are the same as those obtained for the sample covariance (Paul, 2007).

Recall that our high-dimensional results assume an asymptotic regime where $p/n \rightarrow \gamma > 0$, as opposed to the classical regime where p is fixed and $n \rightarrow \infty$. The case of fixed p corresponds to $\gamma = 0$ and the spectral properties of the sample correlation matrix are well understood; see, for example, Girshick (1939), Konishi (1979), Fang and Krishnaiah (1982), Schott (1991), Kollo and Neudecker (1993), and Boik (2003). When $\gamma = 0$, the function $\rho(\ell; \gamma)$ reduces to the identity. Indeed, for fixed p , there is no high-dimensional component η in Model M, and hence no biasing effect on $\rho(\ell; \gamma)$ that occurs when $\gamma > 0$. In particular, for fixed p there is no counterpart to our Theorem 3.

To summarize, in comparison to the high-dimensional ($p/n \rightarrow \gamma > 0$) sample covariance setting, our results for the spiked eigenvalues and eigenvectors of sample correlation matrices confirm that the first-order asymptotic behavior is indeed equivalent to that of sample covariance matrices, in agreement with previous results and observations (El Karoui, 2009, Mestre and Vallet, 2017). While the eigenvalue limits in Theorem 1 and Theorem 3 follow as a straightforward consequence of El Karoui (2009), the eigenvector results of Theorem 2-(i) and Theorem 3-(ii) do not. In contrast to the first-order equivalences, important differences arise in the fluctuations of both the eigenvalues and eigenvectors, as shown by the asymptotic distributions of Theorem 1-(ii) and Theorem 2-(ii).

We illustrate these differences with a simple example having covariance $\Gamma = (1 - r)I_m + r\mathbf{1}_m\mathbf{1}_m^T$, where $r \in [0, 1]$; that is, a model with unit variances and constant correlation r across all components. Moreover, ξ is assumed to be Gaussian for simplicity. In this setting, $L = \text{diag}(\ell_1, 1 - r, \dots, 1 - r)$, where $\ell_1 = 1 + r(m - 1)$ is supercritical iff $r > \sqrt{\gamma}/(m - 1)$. Consider the largest sample eigenvalue $\hat{\ell}_1$ in such a supercritical case. From Corollary 1, the asymptotic variances for the sample covariance and the sample correlation can be computed, yielding

$$\sigma_1^2 = 2\ell_1^2 \dot{\rho}_1, \quad \tilde{\sigma}_1^2 = \sigma_1^2(1 - \dot{\rho}_1 \Delta),$$

respectively, with $\Delta = 2\ell_1 \text{tr} P_D^4 - \text{tr}(P_D \Gamma P_D)^2$, and where

$$P_D \triangleq P_{D,1} = m^{-1/2} I_m, \quad \dot{\rho}_1 = 1 - \frac{\gamma}{r^2(m-1)^2}.$$

Figure 2(a) plots these asymptotic variances versus r for various (γ, m) . Indeed, the variance (fluctuation) for the sample correlation is consistently smaller than for the sample covariance. The difference is striking, becoming extremely large as $r \nearrow 1$. Similar trends are observed for various choices of m and γ , being more pronounced for higher m , while not much affected by varying γ . This may be understood from the fact that, after writing $\dot{\rho}_1 = r(2-r) + (1-r)^2 m^{-1} = 1 - (1-r)^2(1-m^{-1})$,

$$\frac{\tilde{\sigma}_1^2}{\sigma_1^2} = 1 - \dot{\rho}_1 \Delta \rightarrow \begin{cases} \frac{\gamma}{(m-1)^2} & \text{as } r \rightarrow 1, m \text{ fixed} \\ (1-r)^2 & \text{as } m \rightarrow \infty, r \text{ fixed.} \end{cases}$$

Turn now to the fluctuations of the leading sample eigenvector, in the same setting as above. Note that, in Corollary 2, for this particular case, one can deduce from $P^T \Gamma P = L$ that

$$\mathcal{X} = m^{-1}(1-r^2)I_m + r^2 e_1 e_1^T, \quad \mathcal{Y} = m^{-1} I_m.$$

Also from Corollary 2, the asymptotic variances for the normalized sample-to-population eigenvector projection $p_2^T a_1$, in the sample covariance and sample correlation cases, are computed as

$$\Sigma_{1,22}^{\text{cov}} = \frac{\ell_1 \ell_2}{(rm)^2 \dot{\rho}_1}, \quad \Sigma_{1,22} = \Sigma_{1,22}^{\text{cov}} - \frac{\zeta}{(rm)^2} \frac{\ell_1 \ell_2 (\ell_1 + \ell_2)}{m},$$

respectively, where $\zeta = 1 - r + \frac{1}{2}(1+r)\left(1 + \frac{1-r}{rm}\right)^{-1}$, and we recall that $\ell_1 = 1 - r + rm$ and $\ell_2 = 1 - r$. These variances are numerically evaluated in Figure 2(b) for the same parameter choices as before and, again, as functions of r . Note, however, that for better visual appreciation, the range of r has been restricted to supercritical values sufficiently above the critical point $\sqrt{\gamma}/(m-1)$, because the variance explodes at that point. The comparative evaluation again shows smaller variances for the sample correlation. The variance reduction here is less visible in the graphs, because both $\Sigma_{1,22}$ and $\Sigma_{1,22}^{\text{cov}}$ vanish as $r \rightarrow 1$. The ratio, however, behaves quite similarly to the variance ratio $\tilde{\sigma}_1^2/\sigma_1^2$:

$$\frac{\Sigma_{1,22}}{\Sigma_{1,22}^{\text{cov}}} = 1 - \zeta \rho_1 \frac{(\ell_1 + \ell_2)}{m} \rightarrow \begin{cases} \frac{\gamma}{(m-1)^2} & \text{as } r \rightarrow 1, m \text{ fixed} \\ (1-r)(1-r/2) & \text{as } m \rightarrow \infty, r \text{ fixed.} \end{cases}$$

We end the discussion of our main results with a few remarks about possible extensions. Our results assume that $\ell_\nu > 1$ is a simple eigenvalue, but extensions for small spikes with $\ell_\nu < 1$ and for spikes with multiplicities should be possible. Analogous results for eigenvalues have been obtained for sample covariance matrices for $\ell_\nu < 1$, including multiplicities greater than one (e.g., see Bai and Yao (2008)), giving reason to expect corresponding results for correlation matrices. Extensions of our results for eigenvalues and eigenvectors of sample correlation matrices for simple $\ell_\nu < 1$ should be fairly straightforward, though the cases $\gamma < 1$, $\gamma = 1$, and $\gamma > 1$ would need separate treatment. Extensions for spikes with multiplicities are also possible, but in this case the eigenvectors are not well defined and one would need to consider subspace projections, requiring non-trivial modifications of our technical arguments.

The remainder of the paper proceeds as follows. First, in Section 3, we introduce key quantities and identities used in the derivations. Section 4 presents necessary asymptotic properties for bilinear forms and matrix quadratic forms with normalized entries, with the corresponding proofs relegated to the Supplementary Material, Section S3. These properties provide a foundation for describing the asymptotic convergence and distribution of eigenvalues and eigenvectors of sample correlation matrices, derived in Sections 5 and 6 respectively.

As already noted, a parallel treatment for the simpler case of covariance matrices is given in a supplementary manuscript [JY]. This aims at a unified exposition of known spectral properties of spiked covariance matrices as a benchmark for the current work, along with additional citations to the literature.

3. Preliminaries

We begin with a block representation and some associated reductions for the sample correlation matrix R . These are well known in the covariance matrix setting. As with the partition of x in Model M, consider

$$X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}, \quad X_1 \in \mathbb{R}^{m \times n}, \quad X_2 \in \mathbb{R}^{p \times n}.$$

Write $S_D = \text{blkdiag}(S_{D1}, S_{D2})$, with S_{D1} containing the sample variances corresponding to ξ , and S_{D2} containing those corresponding to η . Define the “normalized” data matrices $\bar{X}_1 = S_{D1}^{-1/2} X_1$ and $\bar{X}_2 = S_{D2}^{-1/2} X_2$, such that

$$R = n^{-1} \begin{bmatrix} \bar{X}_1 \bar{X}_1^T & \bar{X}_1 \bar{X}_2^T \\ \bar{X}_2 \bar{X}_1^T & \bar{X}_2 \bar{X}_2^T \end{bmatrix} = \begin{bmatrix} R_{11} & R_{12} \\ R_{21} & R_{22} \end{bmatrix}; \quad \hat{\mathbf{p}}_\nu = \begin{bmatrix} \hat{p}_\nu \\ \hat{u}_\nu \end{bmatrix}.$$

This partitioning of the eigenvector equation $R\hat{\mathbf{p}}_v = \hat{\ell}_v \hat{\mathbf{p}}_v$, along with $\hat{\mathbf{p}}_v = \begin{bmatrix} \hat{p}_v^T \\ \hat{v}_v^T \end{bmatrix}^T$, yields

$$R_{11}\hat{p}_v + R_{12}\hat{v}_v = \hat{\ell}_v \hat{p}_v$$

$$R_{21}\hat{p}_v + R_{22}\hat{v}_v = \hat{\ell}_v \hat{v}_v.$$

From the second equation, $\hat{v}_v = (\hat{\ell}_v I_p - R_{22})^{-1} R_{21} \hat{p}_v$. Substituting this into the first equation yields

$$K(\hat{\ell}_v) \hat{p}_v = \hat{\ell}_v \hat{p}_v, \quad \text{with} \quad K(t) = R_{11} + R_{12}(tI_p - R_{22})^{-1} R_{21}.$$

Thus, $\hat{\ell}_v$ is an eigenvalue of $K(\hat{\ell}_v)$, with associated eigenvector \hat{p}_v ; this is central to our derivations. Note that $K(\hat{\ell}_v)$ is well defined if $\hat{\ell}_v$ is well separated from the eigenvalues of R_{22} ; Section 5.1 shows that this occurs with probability one for all large n when ℓ_v is supercritical. Furthermore, the normalization condition, $\hat{p}_v^T \hat{p}_v + \hat{v}_v^T \hat{v}_v = 1$ yields

$$\hat{p}_v^T (I_m + Q_v) \hat{p}_v = 1, \quad Q_v = R_{12}(\hat{\ell}_v I_p - R_{22})^{-2} R_{21}.$$

Phrased in terms of the signal-space normalized eigenvector $a_v = \hat{p}_v / \|\hat{p}_v\|$, we have

$$K(\hat{\ell}_v) a_v = \hat{\ell}_v a_v, \quad a_v^T (I_m + Q_v) a_v = \|\hat{p}_v\|^{-2}. \tag{3.10}$$

Note also that the sample-to-population inner product can be rewritten as

$$\langle \hat{\mathbf{p}}_v, \mathbf{p}_v \rangle = \langle \hat{p}_v, p_v \rangle = \|\hat{p}_v\| \langle a_v, p_v \rangle. \tag{3.11}$$

In the derivation of our CLT results, we use an eigenvector perturbation formula with quadratic error bound given in [JY, Lemma 13], itself a modification of the arguments in Paul (2007). This yields the key expansion

$$a_v - p_v = -\mathcal{R}_{vn} D_v p_v + r_v, \tag{3.12}$$

where

$$\mathcal{R}_{vn} = \frac{\ell_v}{\rho_{vn}} \sum_{k \neq v}^m (\ell_k - \ell_v)^{-1} p_k p_k^T, \quad D_v = K(\hat{\ell}_v) - (\rho_{vn} / \ell_v) \Gamma, \quad \|r_v\| = O(\|D_v\|^2).$$

The derivations of our eigenvalue and eigenvector results, presented in Sections 5 and 6 respectively, take (3.10), (3.11) and (3.12) as points of departure, and rely on asymptotic

properties of the key objects $K(\hat{\ell}_v)$ and Q_v . In particular, $K(t)$ can be expressed as the random matrix quadratic form

$$K(t) = n^{-1} \bar{X}_1 B_n(t) \bar{X}_1^T, \tag{3.13}$$

where, using the Woodbury identity,

$$\begin{aligned} B_n(t) &= I_n + n^{-1} \bar{X}_2^T (tI_p - R_{22})^{-1} \bar{X}_2 \\ &= t(I_n - n^{-1} \bar{X}_2^T \bar{X}_2)^{-1}. \end{aligned}$$

Thus, our key objects are random quadratic forms involving the normalized data matrices \bar{X}_1 and \bar{X}_2 . The asymptotic properties of these forms are foundational to our results, and are presented next.

4. Quadratic forms with normalized entries

In this section, we establish the first-order (deterministic) convergence and a CLT for matrix quadratic forms of the type $n^{-1} \bar{X}_1 B_n \bar{X}_1^T$, where B_n is a matrix with bounded spectral norm. While being essential to our purposes, some of the technical results may be of independent interest; thus, we first present the general results, and then apply these in the context of Model M.

4.1 First-order convergence

To establish the first-order convergence, we first require some results on bilinear forms involving correlated random vectors of unit length. A main technical result (see Supplementary Material, S3.1) is the following:

Lemma 1—Let B be an $n \times n$ nonrandom symmetric matrix, and let $x, y \in \mathbb{R}^n$ be random vectors of i.i.d. entries with mean zero, variance one, $\mathbb{E}|x_i|^4, \mathbb{E}|y_i|^4 \leq v_4$, and $\mathbb{E}[x_i y_i] = \rho$. Let $\bar{x} = \sqrt{nx}/\|x\|$ and $\bar{y} = \sqrt{ny}/\|y\|$. Then, for any $s \geq 1$,

$$\mathbb{E} \left| n^{-1} \bar{x}^T B \bar{y} - \rho n^{-1} \text{tr} B \right|^s \leq \mathcal{C}_s \left[n^{-s} \left(v_2 s \text{tr} B^s + (v_4 \text{tr} B^2)^{s/2} \right) + \|B\|^s \left(n^{-s/2} v_4^{s/2} + n^{-s} + v_2 s \right) \right],$$

where \mathcal{C}_s is a constant depending only on s .

This is a generalization of Gao et al. (2017, Lemma 5), which established a corresponding bound for normalized quadratic forms. Lemma 1 leads to the following first-order convergence result:

Corollary 3—Let x, y be random vectors of i.i.d. entries with mean zero, variance one, $\mathbb{E}|x_i|^{4+\delta}, \mathbb{E}|y_i|^{4+\delta} < \infty$ for some $\delta > 0$, and $\mathbb{E}[x_i y_i] = \rho$. Define $\bar{x} = \sqrt{nx}/\|x\|$ and $\bar{y} = \sqrt{ny}/\|y\|$, and let B_n be a sequence of $n \times n$ symmetric matrices, with $\|B_n\|$ bounded. Then,

$$n^{-1} \bar{x}^T B_n \bar{y} - n^{-1} \rho \text{tr} B_n \xrightarrow{\text{a.s.}} 0.$$

Proof. Because the $(4 + \delta)$ th moment and $\|B_n\|$ are bounded, from Lemma 1,

$$\mathbb{E} \left| n^{-1} \bar{x}^T B_n \bar{y} - n^{-1} \rho \text{tr} B_n \right|^{2 + \delta/2} \leq O(n^{-(1 + \delta/4)}).$$

The convergence then follows from Markov’s inequality and the Borel–Cantelli lemma. \square

We now apply this to our Model M with *random* matrices $B_n(\bar{X}_2)$, independent of \bar{X}_1 :

Lemma 2—Assume Model M, and suppose that $B_n = B_n(\bar{X}_2)$ is a sequence of random symmetric matrices, for which $\|B_n\|$ is $O_{\text{a.s.}}(1)$. Then,

$$n^{-1} \bar{X}_1 B_n(\bar{X}_2) \bar{X}_1^T - n^{-1} \text{tr} B_n(\bar{X}_2) \Gamma \xrightarrow{\text{a.s.}} 0.$$

Proof. This follows from Fubini’s theorem. Specifically, one may use the arguments in the proof of [JY, Lemma 5], applying Corollary 3, and noting that \bar{X}_1 is independent of $B_n(\bar{X}_2)$. \square

4.2 Central Limit Theorem

To establish our main matrix quadratic-form CLT result, we first derive a CLT for scalar bilinear forms involving normalized random vectors. To this end, we must introduce some further notation. Consider zero-mean random vectors $(x, y) \in \mathbb{R}^M \times \mathbb{R}^M$, with

$$\text{Cov} \begin{pmatrix} x \\ y \end{pmatrix} = C = \begin{pmatrix} C^{xx} & C^{xy} \\ C^{yx} & C^{yy} \end{pmatrix},$$

where $C_{ll'}^{xy} = \mathbb{E}[x_l y_{l'}]$. Assume $C_{ll}^{xx} = C_{ll}^{yy} = 1$; that is, all components of the x and y vectors have unit variance and $\rho_l = C_{ll}^{xy} = \mathbb{E}[x_l y_l]$. We first introduce notation for some quadratic functions of x_l, y_l . Let $z, w \in \mathbb{R}^M$, with

$$z_l = x_l y_l, \quad w_l = \rho_l (x_l^2 + y_l^2)/2, \quad C^{zz} = \text{Cov}(z), \quad C^{wz} = \text{Cov}(z, w), \text{ etc.}$$

Let $X = (x_{ij})_{M \times n}$ and $Y = (y_{ij})_{M \times n}$ be data matrices based on n i.i.d. observations of (x, y) , and define the “normalized” data matrices $\bar{X} = \hat{\Sigma}_x^{-1/2} X$ and $\bar{Y} = \hat{\Sigma}_y^{-1/2} Y$, where

$$\hat{\Sigma}_x = \text{diag}(\hat{\sigma}_{x_1}^2, \dots, \hat{\sigma}_{x_M}^2), \quad \hat{\Sigma}_y = \text{diag}(\hat{\sigma}_{y_1}^2, \dots, \hat{\sigma}_{y_M}^2), \quad \text{and} \quad \hat{\sigma}_{x_l}^2 = n^{-1} \sum_{i=1}^n x_{il}^2, \hat{\sigma}_{y_l}^2 = n^{-1} \sum_{i=1}^n y_{il}^2.$$

Then, we use the following notation for the rows \bar{x}_l^T and \bar{y}_l^T of the normalized data matrices

$$\bar{X} = (\bar{x}_l)_{M \times n} = \begin{bmatrix} \bar{x}_1^T \\ \vdots \\ \bar{x}_M^T \end{bmatrix}, \quad \bar{Y} = (\bar{y}_l)_{M \times n} = \begin{bmatrix} \bar{y}_1^T \\ \vdots \\ \bar{y}_M^T \end{bmatrix}.$$

With this setup, we have the following result, proved in the Supplementary Material, S3.2:

Proposition 1—Let $B_n = (b_{n,ij})$ be random symmetric $n \times n$ matrices, independent of X, Y , such that for some finite $\beta, \|B_n\| \leq \beta$ for all n , and

$$n^{-1} \sum_{i=1}^n b_{n,ii}^2 \xrightarrow{p} \omega, \quad n^{-1} \text{tr} B_n^2 \xrightarrow{p} \theta, \quad (n^{-1} \text{tr} B_n)^2 \xrightarrow{p} \phi,$$

all finite. In addition, define $Z_n \in \mathbb{R}^M$, with components

$$Z_{n,l} = n^{-1/2} [\bar{x}_l^T B_n \bar{y}_l - \rho_l \text{tr} B_n].$$

Then, $Z_n \xrightarrow{\mathcal{D}} N_M(0, D)$, with

$$D = (\theta - \omega)J + \omega K_1 + \phi K_2 = \theta J + \omega K + \phi K_2, \tag{4.14}$$

where $K = K_1 - J$ and J, K_1, K_2 are matrices defined by

$$\begin{aligned} J &= C^{xy} \circ C^{yx} + C^{xx} \circ C^{yy} \\ K_1 &= C^{zz} \\ K_2 &= C^{ww} - C^{wz} - C^{zw}. \end{aligned} \tag{4.15}$$

The entries of K are fourth-order cumulants of x and y :

$$K_{ll'} = \mathbb{E}(x_l y_l x_{l'} y_{l'}) - \mathbb{E}(x_l y_l) \mathbb{E}(x_{l'} y_{l'}) - \mathbb{E}(x_l y_{l'}) \mathbb{E}(y_l x_{l'}) - \mathbb{E}(x_l x_{l'}) \mathbb{E}(y_l y_{l'}). \tag{4.16}$$

Hence, K vanishes if x, y are Gaussian.

The corresponding result with unnormalized vectors is established in [JY Theorem 10]. The terms $\theta J + \omega K$ appear in that case, and the additional term ϕK_2 reflects the normalization in \bar{x}_l and \bar{y}_l . As in [JY], the proof is based on the martingale CLT, rather than the moment method used in Bai and Yao (2008), which stated a similar result for quadratic forms involving unnormalized random vectors.

While potentially of independent interest, Proposition 1 is important for our purposes through its application to Model M.

Proposition 2—Assume Model M, and consider B_n as in Proposition 1. Then,

$$W_n = n^{-1/2} [\bar{X}_1 B_n \bar{X}_1^T - (\text{tr} B_n) \Gamma] \xrightarrow{\mathcal{D}} W,$$

where W is a symmetric $m \times m$ Gaussian matrix with entries W_{ij} , mean zero, and covariances given by

$$\text{Cov}[W_{ij}, W_{i'j'}] = \theta(\kappa_{ij'}\kappa_{ji'} + \kappa_{ii'}\kappa_{jj'}) + \omega\kappa_{ijj'j'} + \phi\check{\kappa}_{ijj'j'}, \tag{4.17}$$

for $i \neq j$ and $i' \neq j'$.

Proof. The result follows from Proposition 1 by turning the matrix quadratic form $\bar{X}_1 B_n \bar{X}_1^T$ into a vector of bilinear forms; see, for example, [JY, Proposition 6] and Bai and Yao (2008, Proposition 3.1). Specifically, use an index l for the $M = m(m+1)/2$ pairs (i, j) , with $1 \leq i < j \leq m$. Build the random vectors (x, y) for Proposition 1 as follows: if $l = (i, j)$, then set $x_l = \xi_i/\sigma_i$ and $y_l = \xi_j/\sigma_j$. In the resulting covariance matrix C for (x, y) , if also $l' = (i', j')$,

$$C_{ll'}^{xy} = \mathbb{E}[\xi_i \xi_{j'}] / (\sigma_i \sigma_{j'}) = \kappa_{ij'}, \quad C_{ll'}^{yx} = \kappa_{ji'}, \quad C_{ll'}^{xx} = \kappa_{ii'}, \quad C_{ll'}^{yy} = \kappa_{jj'}$$

and, in particular, $\rho_l = C_{ll}^{xy} = \kappa_{ij}$ and $\rho_{l'} = \kappa_{i'j'}$, whereas $C_{ll}^{xx} = C_{ll}^{yy} = 1$. Component $W_{n,ij}$ corresponds to component Z_l in Proposition 1. Thus, we conclude that $W_n \xrightarrow{\mathcal{D}} W$, where W is a Gaussian matrix with zero mean and $\text{Cov}(W_{ij}, W_{i'j'}) = D_{ll'}$, given by Proposition 1. It remains to interpret the quantities in (4.14) in terms of Model M. Substituting $x_l = \bar{\xi}_i$ and $y_l = \bar{\xi}_j$ into (4.16) and chasing definitions, we obtain $J_{ll'} = \kappa_{ij'}\kappa_{ji'} + \kappa_{ii'}\kappa_{jj'}$ and $K_{ll'} = \kappa_{ijj'j'}$. Observing that $z_l = x_l y_l = \chi_{ij}$ and $w_l = \rho_l(x_l^2 + y_l^2)/2 = \psi_{ij}$, we similarly find that $K_{2, ll'} = \check{\kappa}_{ijj'j'}$. \square

5. Proofs of the eigenvalue results

In this section, we derive the main eigenvalue results, presented in Theorem 1 and Theorem 3-(i).

5.1 Preliminaries

Convergence properties of the eigenvalues of R_{22} —It is well known that the empirical spectral density (ESD) of S_{22} converges weakly a.s. to the Marchenko–Pastur (MP) law F_γ , and that the extreme non-trivial eigenvalues converge to the edges of the support of F_γ . For the sample correlation case, Jiang (2004b) shows that the same is true for R_{22} . That is, the empirical distribution of the eigenvalues $\mu_1 \dots \mu_p$ of the “noise” correlation matrix $R_{22} = n^{-1} \bar{X}_2 \bar{X}_2^T$ converges weakly a.s. to the MP law F_γ , supported on $[a_\gamma, b_\gamma] = [(1 - \sqrt{\gamma})^2, (1 + \sqrt{\gamma})^2]$, if $\gamma \leq 1$, and on $\{0\} \cup [a_\gamma, b_\gamma]$ otherwise. Furthermore, the ESD of the $n \times n$ companion matrix $C_n = n^{-1} \bar{X}_2^T \bar{X}_2$, denoted by F_n , converges weakly a.s. to

the “companion MP law” $F_\gamma = (1 - \gamma)\mathbf{1}_{[0,\infty)} + \gamma F_\gamma$, where $\mathbf{1}_A$ denotes the indicator function on set A .

In addition, Jiang (2004b) shows that

$$\mu_1 \xrightarrow{\text{a.s.}} b_\gamma \quad \text{and} \quad \mu_{p \wedge n} \xrightarrow{\text{a.s.}} a_\gamma. \tag{5.18}$$

Based on these results, if $f_n \rightarrow f$ uniformly as continuous functions on the closure \mathcal{S} of a bounded neighborhood of the support of F_γ , then:

$$\int f_n(x)F_n(dx) \xrightarrow{\text{a.s.}} \int f(x)F_\gamma(dx). \tag{5.19}$$

If $\text{supp}(F_n)$ is not contained in \mathcal{S} , then the left side integral may not be defined. However, such an event occurs for at most finitely many n with probability one.

Almost sure limit of $\hat{\ell}_v$ —The statements in Theorem 1-(i) and Theorem 3-(i) follow easily from known results. Specifically, denote the v th eigenvalue of the sample covariance S by $\hat{\lambda}_v$. The almost sure limits

$$\hat{\lambda}_v \xrightarrow{\text{a.s.}} \begin{cases} \rho_v, & \ell_v > 1 + \sqrt{\gamma} \\ (1 + \sqrt{\gamma})^2, & 1 < \ell_v \leq 1 + \sqrt{\gamma} \end{cases} \tag{5.20}$$

were established in Baik and Silverstein (2006). From the proof of El Karoui (2009, Lemma 1),

$$\max_{i=1, \dots, m} |\hat{\lambda}_i - \ell_i| \xrightarrow{\text{a.s.}} 0.$$

Therefore, the same almost sure limits as (5.20) hold for $\hat{\ell}_v$.

High-probability events $J_{n\epsilon}$, $J_{n\epsilon 1}$ —When necessary, we may confine attention to the event $J_{n\epsilon} = \{\hat{\ell}_v > \min(\rho_v, \rho_{vn}) - \epsilon, \mu_1 \leq b_\gamma + \epsilon\}$ or $J_{n\epsilon 1} = \{\mu_1 \leq b_\gamma + \epsilon\}$, with $\epsilon > 0$ chosen such that $\rho_v - b_\gamma \geq 3\epsilon$, because from (2.5) (proven above) and (5.18), these events occur with probability one for all large n .

Asymptotic expansion of $K(\hat{\ell}_v)$ —We establish an asymptotic stochastic expansion for the quadratic form $K(\hat{\ell}_v)$. Specifically, using the decomposition

$$K(\hat{\ell}_v) = K(\rho_{vn}) + [K(\hat{\ell}_v) - K(\rho_{vn})], \tag{5.21}$$

we show that

$$K(\rho_{vn}) \xrightarrow{\text{a.s.}} -\rho_v m(\rho_v; \gamma) \Gamma = (\rho_v / \ell_v) \Gamma \tag{5.22}$$

and

$$K(\widehat{\ell}_v) - K(\rho_{vn}) = -(\widehat{\ell}_v - \rho_{vn})[c(\rho_v)\Gamma + o_{a.s.}(1)], \tag{5.23}$$

where, for $t \notin \text{supp}(F_\gamma)$,

$$m(t; \gamma) = \int (x - t)^{-1} F_\gamma(dx), \quad c(t) = \int x(t - x)^{-2} F_\gamma(dx).$$

Here, m is the Stieltjes transform of the companion distribution F_γ .

In establishing (5.22), start by taking sufficiently large n such that $|\rho_{vn} - \rho_v| < \epsilon$, with ϵ defined as above. For such n , on $J_{n \in I}$, we have

$$\|B_n(\rho_{vn})\| \leq \frac{\rho_v + \epsilon}{\epsilon}.$$

Because $J_{n \in I}$ holds with probability one for all large n , $\|B_n(\rho_{vn})\| = O_{a.s.}(1)$ and, therefore, it follows from Lemma 2 that

$$K(\rho_{vn}) - n^{-1} \text{tr} B_n(\rho_{vn}) \Gamma \xrightarrow{a.s.} 0.$$

In addition, (5.19) yields

$$n^{-1} \text{tr} B_n(\rho_{vn}) = \int \rho_{vn}(\rho_{vn} - x)^{-1} F_n(dx) \xrightarrow{a.s.} \int \rho_v(\rho_v - x)^{-1} F_\gamma(dx) = -\rho_v m(\rho_v; \gamma).$$

Explicit evaluation gives $m(\rho_v; \gamma) = -1/\ell_v$, [JY, Appendix A], and (5.22) follows.

To establish (5.23), we start by recalling that $C_n = n^{-1} \bar{X}_2^T \bar{X}_2$, and introduce the resolvent notation $Z(t) = (tI_n - C_n)^{-1}$, such that $B_n(t) = tZ(t)$ and $K(t) = n^{-1} \bar{X}_1^T tZ(t) \bar{X}_1^T$. From the resolvent identity, that is, $A^{-1} - B^{-1} = A^{-1}(B - A)B^{-1}$ for square invertible A and B , and noting that $tZ(t) = C_n Z(t) + I$ from the Woodbury identity, we have, for $t_1, t_2 > b_\gamma$,

$$t_1 Z(t_1) - t_2 Z(t_2) = -(t_1 - t_2) C_n Z(t_1) Z(t_2)$$

and, therefore,

$$K(\widehat{\ell}_v) - K(\rho_{vn}) = -(\widehat{\ell}_v - \rho_{vn}) n^{-1} \bar{X}_1^T C_n Z(\widehat{\ell}_v) Z(\rho_{vn}) \bar{X}_1^T.$$

Moreover, again by the resolvent identity, $Z(\widehat{\ell}_v) = Z(\rho_{vn}) - (\widehat{\ell}_v - \rho_{vn}) Z(\widehat{\ell}_v) Z(\rho_{vn})$, which yields

$$\begin{aligned}
 K(\widehat{\ell}_v) - K(\rho_{vn}) &= -(\widehat{\ell}_v - \rho_{vn})n^{-1}\bar{X}_1 B_{n1}(\rho_{vn}, \rho_{vn})\bar{X}_1^T \\
 &+ (\widehat{\ell}_v - \rho_{vn})^2 n^{-1}\bar{X}_1 B_{n2}(\widehat{\ell}_v, \rho_{vn})\bar{X}_1^T,
 \end{aligned} \tag{5.24}$$

with $B_{nr}(t_1, t_2)$ defined as

$$B_{nr}(t_1, t_2) = C_n Z(t_1) Z^r(t_2). \tag{5.25}$$

We now characterize the first-order behavior of the two matrix quadratic forms in (5.24). For the first, we simply mirror the arguments of the proof of (5.22) to obtain

$$n^{-1}\bar{X}_1 B_{n1}(\rho_{vn}, \rho_{vn})\bar{X}_1^T \xrightarrow{\text{a.s.}} c(\rho_v)\Gamma.$$

For the second, we again apply similar reasoning, operating on the event $J_{n\epsilon}$. Specifically, it is easy to establish that on $J_{n\epsilon}$, and for n sufficiently large that $|\rho_{vn} - \rho_v| < \epsilon$, $\|B_{n2}(\widehat{\ell}_v, \rho_{vn})\|$ is bounded. Hence, $\|B_{n2}(\widehat{\ell}_v, \rho_{vn})\| = O_{\text{a.s.}}(1)$, and it follows from Lemma 2 and (5.19) that

$$n^{-1}\bar{X}_1 B_{n2}(\widehat{\ell}_v, \rho_{vn})\bar{X}_1^T = O_{\text{a.s.}}(1).$$

The expansion in (5.23) is obtained by combining the latter two equations with (5.24).

CLT of $K(\rho_{vn})$ —We now specialize Proposition 2 for the matrix quadratic form $K(\rho_{vn})$.

Proposition 3—Assume Model M, and define ρ_{vn} by (1.1) and $K(\rho_{vn})$ by (3.13). Then,

$$W_n(\rho_{vn}) = \sqrt{n}[K(\rho_{vn}) - n^{-1}\text{tr}B_n(\rho_{vn})\Gamma] \xrightarrow{\mathcal{D}} W^v,$$

which is a symmetric Gaussian random matrix with entries W_{ij}^v , mean zero, and covariances given by

$$\text{Cov}[W_{ij}^v, W_{i'j'}^v] = \frac{\rho_v^2}{\ell_v^2 \rho_v} (\kappa_{ij'} \kappa_{ji'} + \kappa_{ii'} \kappa_{jj'}) + \frac{\rho_v^2}{\ell_v^2} (\kappa_{iji'j'} + \check{\kappa}_{iji'j'}), \tag{5.26}$$

where ρ_v and $\dot{\rho}_v$ are defined in (1.1), and the terms in parentheses are defined in (1.2) and (1.4).

Proof. Recall that $J_{n\epsilon} = \{\mu_1 > b_\gamma + \epsilon\}$, and consider sufficiently large n such that $\rho_{vn} > \rho_v - \epsilon$. Then, we may apply Proposition 2 with $B_n = B_n(\rho_{vn})1_{J_{n\epsilon}}$, which is independent of \bar{X}_1 , and for which $\|B_n\|$ is bounded. Specifically, the result follows by applying Proposition 2 to $W_n(\rho_{vn})1_{J_{n\epsilon}}$, along with the fact that $1_{J_{n\epsilon}} \xrightarrow{\text{a.s.}} 1$, and particularizing ω , θ , and ϕ in (4.17). These quantities, denoted respectively by ω_v , θ_v , and ϕ_v , can be computed as in [JY, Appendix A], yielding

$$\omega_v = \phi_v = \frac{(\ell_v - 1 + \gamma)^2}{(\ell_v - 1)^2} = \frac{\rho_v^2}{\ell_v^2}, \quad \theta_v = \frac{(\ell_v - 1 + \gamma)^2}{(\ell_v - 1)^2 - \gamma} = \frac{\omega_v}{\rho_v}.$$

Tightness properties—Lastly, we establish some tightness properties essential to the derivation of our second-order results.

We first establish a refinement of (5.22). Define $K_0(\rho; \gamma) := -\rho m(\rho; \gamma)\Gamma$, such that (5.22) is rewritten as $K(\rho_{vn}) \xrightarrow{\text{a.s.}} K_0(\rho_v; \gamma)$. Set $g_\rho(x) = \rho(\rho - x)^{-1}$, and write

$$\text{tr} B_n(\rho) = \sum_{i=1}^n \rho(\rho - \mu_i)^{-1} = \sum_{i=1}^n g_\rho(\mu_i).$$

In addition, introducing

$$G_n(g) := \sum_{i=1}^n g(\mu_i) - n \int g(x) F_{\gamma_n}(dx),$$

we have

$$\begin{aligned} K(\rho) - K_0(\rho; \gamma_n) &= K(\rho) - n^{-1} \text{tr} B_n(\rho)\Gamma \\ &+ \rho n^{-1} \left[\sum_{i=1}^n (\rho - \mu_i)^{-1} - n \int (\rho - x)^{-1} F_{\gamma_n}(dx) \right] \Gamma \end{aligned} \tag{5.27}$$

$$= n^{-1/2} W_n(\rho) + n^{-1} G_n(g_\rho)\Gamma.$$

Lemma 3—Assume that Model M holds, and that $\ell_v > 1 + \sqrt{\gamma}$ is simple. For some $b > \rho_1$, let I denote the interval $[b_\gamma + 3\epsilon, b]$. Then,

$$\{G_n(g_\rho), \rho \in I\} \text{ is uniformly tight,} \tag{5.28}$$

$$\{n^{1/2}[K(\rho) - K_0(\rho; \gamma_n)], \rho \in I\} \text{ is uniformly tight,} \tag{5.29}$$

$$\widehat{\ell}_v - \rho_{vn} = O_p(n^{-1/2}), \tag{5.30}$$

$$a_v - p_v = O_p(n^{-1/2}). \tag{5.31}$$

Proof. The proofs of (5.28)–(5.30) appear in the Supplementary Material, S2. We show (5.31) using the expansion $a_v - p_v = -\mathcal{R}_{vn} D_v p_v + r_v$, given in (3.12), from which we recall $\|r_v\| = O(\|D_v\|^2)$ and note that $\|\mathcal{R}_{vn}\| \leq C$ and $D_v = K(\widehat{\ell}_v) - K_0(\rho_{vn}; \gamma_n)$. We then have $a_v - p_v = O_p(\|D_v\| + \|D_v\|^2)$. Furthermore, from

$$\|D_v\| \leq \|K(\hat{\ell}_v) - K(\rho_{vn})\| + \|K(\rho_{vn}) - K_0(\rho_{vn}; \gamma_n)\|,$$

the first term is $O_p(n^{-1/2})$ by (5.23) and (5.30), as is the second term by (5.29). Hence,

$$\|D_v\| = O_p(n^{-1/2}), \tag{5.32}$$

and the proof is completed. \square

5.2 Eigenvalue fluctuations (Theorem 1-(ii))

The proof of Theorem 1-(ii) relies on the key expansion

$$\sqrt{n}(\hat{\ell}_v - \rho_{vn})[1 + c(\rho_v)\ell_v + o_p(1)] = p_v^T W_n(\rho_{vn})p_v + o_p(1), \tag{5.33}$$

which is obtained by combining the vector equations $K(\hat{\ell}_v)a_v = \hat{\ell}_v a_v$ and $K_0(\rho_{vn}; \gamma_n)p_v = \rho_{vn}p_v$ with expansions (5.24) for $K(\hat{\ell}_v) - K(\rho_{vn})$ and (5.27) for $K(\rho_{vn}) - K_0(\rho_{vn}; \gamma_n)$.

Specifically, we first use $[K(\hat{\ell}_v) - \hat{\ell}_v I_m]a_v = 0$ to obtain

$$p_v^T [K(\hat{\ell}_v) - \hat{\ell}_v I_m]p_v = (a_v - p_v)^T [K(\hat{\ell}_v) - \hat{\ell}_v I_m](a_v - p_v) = O_p(n^{-1}), \tag{5.34}$$

because $\|K(\hat{\ell}_v) - \hat{\ell}_v I_m\| = O_p(1)$ from (5.21)–(5.23) and (2.5), and $a_v - p_v = O_p(n^{-1/2})$ from Lemma 3. In addition, because $[K_0(\rho_{vn}; \gamma_n) - \rho_{vn}I_m]p_v = 0$, it follows that

$$\begin{aligned} p_v^T [K(\hat{\ell}_v) - \hat{\ell}_v I_m]p_v &= p_v^T [K(\hat{\ell}_v) - K_0(\rho_{vn}; \gamma_n) - (\hat{\ell}_v - \rho_{vn})I_m]p_v \\ &= p_v^T [K(\hat{\ell}_v) - K(\rho_{vn}) - (\hat{\ell}_v - \rho_{vn})I_m]p_v \\ &+ p_v^T [K(\rho_{vn}) - K_0(\rho_{vn}; \gamma_n)]p_v \\ &= -(\hat{\ell}_v - \rho_{vn})[1 + c(\rho_v)\ell_v + o_p(1)] \\ &+ n^{-1/2} p_v^T W_n(\rho_{vn})p_v + o_p(n^{-1/2}), \end{aligned} \tag{5.35}$$

where the last equality follows from (5.23), (5.27), and (5.28). Combining (5.34) and (5.35) yields (5.33).

The asymptotic normality of $\sqrt{n}(\hat{\ell}_v - \rho_{vn})$ now follows from Proposition 3, with asymptotic variance

$$\tilde{\sigma}_v^2 = [1 + c(\rho_v)\ell_v]^{-2} \text{Var}[p_v^T W^v p_v] = (\dot{\rho}_v \ell_v / \rho_v)^2 \sum_{i, j, i', j'} \mathcal{P}_{ij i' j'}^v \text{Cov}[W_{ij}^v, W_{i' j'}^v],$$

where W^v is the $m \times m$ symmetric Gaussian random matrix defined in Proposition 3, with covariance $\text{Cov}[W_{ij}^v, W_{i' j'}^v]$ given by (5.26). Using this in the developed expression for the variance above leads to

$$\tilde{\sigma}_v^2 = \dot{\rho}_v \sum_{i, j, i', j'} \mathcal{P}_{ij i' j'}^v (\kappa_{ij' j' i'} + \kappa_{ii' j' j'}) + \dot{\rho}_v^2 [\mathcal{P}^v, \kappa + \check{\kappa}]. \tag{5.36}$$

By symmetry and the eigen equation $(\Gamma p_v)_i = \sum_j \kappa_{ij} p_{v,j} = \ell_v p_{v,i}$, we have

$$\sum_{i,j,i',j'} \mathcal{P}_{ij i' j'}^v \kappa_{ii'} \kappa_{jj'} = \sum_{i,j,i',j'} \mathcal{P}_{ij i' j'}^v \kappa_{ij'} \kappa_{ji'} = \sum_{i,j} p_{v,i} p_{v,j} (\Gamma p_v)_i (\Gamma p_v)_j = \ell_v^2 \sum_{i,j} (p_{v,i} p_{v,j})^2 = \ell_v^2.$$

Therefore, the first sum in (5.36) reduces to $2\hat{\rho}_v \ell_v^2$, yielding formula (2.6) of Theorem 1.

6. Proofs of the eigenvector results

We now derive the main eigenvector results, presented in Theorem 2 and Theorem 3-(ii).

6.1 Eigenvector inconsistency (Theorem 2-(i))

The convergence result of Theorem 2-(i) follows from two facts: $a_v \xrightarrow{\text{a.s.}} p_v$ and $Q_v \xrightarrow{\text{a.s.}} c(\rho_v)\Gamma$, which are shown below. Once these facts are established, from (3.10),

$$\|\hat{p}_v\|^{-2} \xrightarrow{\text{a.s.}} p_v^T (I_m + c(\rho_v)\Gamma) p_v = 1 + c(\rho_v)\ell_v = \frac{\rho_v}{\ell_v \hat{\rho}_v},$$

which leads to

$$\text{a.s.} \lim \langle \hat{p}_v, p_v \rangle^2 = \text{a.s.} \lim \langle \hat{p}_v, p_v \rangle^2 = \text{a.s.} \lim \|\hat{p}_v\|^2 = \frac{\ell_v \hat{\rho}_v}{\rho_v}.$$

Proof of $a_v \xrightarrow{\text{a.s.}} p_v$ —This is a direct consequence of (3.12) and

$$D_v = K(\rho_{vn}) - (\rho_{vn} \ell_v)\Gamma + K(\hat{\ell}_v) - K(\rho_{vn}) \xrightarrow{\text{a.s.}} 0,$$

which follows from (5.22), (5.23), and the fact that $\hat{\ell}_v - \rho_{vn} \xrightarrow{\text{a.s.}} 0$, given in (2.5).

Proof of $Q_v \xrightarrow{\text{a.s.}} c(\rho_v)\Gamma$ —With $\check{Z}(t) = (tI_p - R_{22})^{-1}$, we have

$$Q_v = R_{12} \check{Z}^2(\rho_v) R_{21} + R_{12} \left[\check{Z}^2(\hat{\ell}_v) - \check{Z}^2(\rho_v) \right] R_{21} \triangleq Q_{v1} + Q_{v2}.$$

Rewrite $Q_{v1} = n^{-1} \bar{X}_1 \check{B}_{n1} \bar{X}_1^T$, with $\check{B}_{n1} = n^{-1} \bar{X}_2^T \check{Z}^2(\rho_v) \bar{X}_2$. On the high-probability event $J_{n\epsilon 1} = \{\mu_1 - b_\gamma + \epsilon\}$, with $\epsilon > 0$ such that $\rho_v - b_\gamma > 2\epsilon$, it is easily established that $\|\check{B}_{n1}\|$ is bounded and, consequently, that $\|\check{B}_{n1}\| = O_{\text{a.s.}}(1)$. Hence, Lemma 2 can be applied to Q_{v1} . Moreover, from (5.19) and noting that

$$n^{-1} \text{tr} \check{B}_{n1} = n^{-1} \text{tr} B_{n1}(\rho_v, \rho_v),$$

with B_{n1} defined in (5.25), we have

$$n^{-1} \text{tr} \check{B}_{n1} \xrightarrow{\text{a.s.}} \int x(\rho_v - x)^{-2} F_{\gamma}(dx) = c(\rho_v).$$

This and Lemma 2 imply that $Q_{v1} \xrightarrow{\text{a.s.}} c(\rho_v)\Gamma$.

It remains to show $Q_{v2} \xrightarrow{\text{a.s.}} 0$. Using a variant of the resolvent identity, that is, $A^{-2} - B^{-2} = -A^{-2}(A^2 - B^2)B^{-2}$ for square invertible A and B , we rewrite

$$Q_{v2} = -2(\hat{\ell}_v - \rho_v)n^{-1}\bar{X}_1\check{B}_{n2}\bar{X}_1^T,$$

with $\check{B}_{n2} = n^{-1}\bar{X}_2^T\check{Z}^2(\hat{\ell}_v)\left[\frac{1}{2}(\hat{\ell}_v + \rho_v)I - R_{22}\right]\check{Z}^2(\rho_v)\bar{X}_2$. Working on the high-probability event $J_{n\epsilon}$, it can be verified that $\|\check{B}_{n2}\| = O_{\text{a.s.}}(1)$. Thus, Lemma 2 together with (5.19) imply that $n^{-1}\bar{X}_1\check{B}_{n2}\bar{X}_1^T = O_{\text{a.s.}}(1)$. Because $\hat{\ell}_v \xrightarrow{\text{a.s.}} \rho_v$, we conclude that $Q_{v2} \xrightarrow{\text{a.s.}} 0$.

6.2 Eigenvector fluctuations (Theorem 2-(ii))

Again, we use the key expansion (3.12). Because $\|r_{\sqrt{n}}\| = O(\|D_{\sqrt{n}}\|^2) = O_p(n^{-1})$ from (5.32), we have

$$\sqrt{n}(a_v - p_v) = -\mathcal{R}_{v\sqrt{n}}\sqrt{n}D_v p_v + o_p(1).$$

Furthermore, using a similar decomposition to the derivation of (5.35),

$$\begin{aligned} \sqrt{n}D_v &= \sqrt{n}[K(\hat{\ell}_v) - K(\rho_{v\sqrt{n}})] + \sqrt{n}[K(\rho_{v\sqrt{n}}) - K_0(\rho_{v\sqrt{n}}, \gamma_n)] \\ &= W_n(\rho_{v\sqrt{n}}) - \sqrt{n}(\hat{\ell}_v - \rho_{v\sqrt{n}})c(\rho_v)\Gamma + o_p(1), \end{aligned}$$

where we use (5.23) and (5.27), along with (5.28) and (5.30) of Lemma 3. Hence, noting that $\mathcal{R}_{v\sqrt{n}}\Gamma p_v = \ell_v \mathcal{R}_{v\sqrt{n}} p_v = 0$ from the definition of $\mathcal{R}_{v\sqrt{n}}$ in (3.12), we have

$$\sqrt{n}(a_v - p_v) = -\mathcal{R}_{v\sqrt{n}}W_n(\rho_{v\sqrt{n}})p_v + o_p(1),$$

or equivalently,

$$\sqrt{n}(P^T a_v - e_v) = -\bar{\mathcal{R}}_{v\sqrt{n}}\bar{W}_n(\rho_{v\sqrt{n}})e_v + o_p(1),$$

where

$$\bar{\mathcal{R}}_{v\sqrt{n}} = \frac{\ell_v}{\rho_{v\sqrt{n}}} \sum_{k \neq v}^m (\ell_k - \ell_v)^{-1} e_k e_k^T, \quad \bar{W}_n(\rho_{v\sqrt{n}}) = P^T W_n(\rho_{v\sqrt{n}})P.$$

The CLT for $P^T a_v$ now follows from Proposition 3. In particular,

$$\sqrt{n}(P^T a_v - e_v) \xrightarrow{\mathcal{D}} \tilde{\mathcal{R}}_v w_v \sim N(0, \Sigma_v),$$

where $\tilde{\mathcal{R}}_v = (\ell_v / \rho_v) \mathcal{D}_v$, recall (2.8), and $w_v = P^T W^\nu p_v$, with W^ν defined in Proposition 3. The covariance matrix $\Sigma_v = \tilde{\mathcal{R}}_v \mathbb{E}[w_v w_v^T] \tilde{\mathcal{R}}_v = \mathcal{D}_v \tilde{\Sigma}_v \mathcal{D}_v$, with $\tilde{\Sigma}_v = (\ell_v / \rho_v)^2 \mathbb{E}[w_v w_v^T]$. The k th component of w_v is given by $w_v(k) = p_k^T W^\nu p_v = \sum_{i,j} p_{k,i} p_{v,j} W_{ij}^\nu$ and, therefore,

$$\tilde{\Sigma}_{v,kl} = \sum_{i,j,i',j'} p_{k,i} p_{v,j} p_{l,i'} p_{v,j'} (\ell_v / \rho_v)^2 \text{Cov}[W_{ij}^\nu, W_{i'j'}^\nu]. \tag{6.37}$$

Theorem 2-(ii) follows after substituting (5.26) for $\text{Cov}[W_{ij}^\nu, W_{i'j'}^\nu]$ and noting that, when $k, l \neq v$,

$$\sum_{i,j,i',j'} p_{k,i} p_{v,j} p_{l,i'} p_{v,j'} (\kappa_{ii' \kappa_{jj'}} + \kappa_{ij' \kappa_{ji'}}) = p_k^T \Gamma p_l \cdot p_v^T \Gamma p_v + p_k^T \Gamma p_v \cdot p_v^T \Gamma p_l = \delta_{kl} \ell_k \ell_v.$$

6.3 Eigenvector inconsistency in the subcritical case (Theorem 3-(ii))

From (3.10) and (3.11), it suffices to show that $a_v^T Q_v a_v \xrightarrow{\text{a.s.}} \infty$ in order for Theorem 3-(ii) to hold. We establish this by showing that $\lambda_{\min}(Q_v) \xrightarrow{\text{a.s.}} \infty$. The approach uses a regularized version of Q_v ,

$$Q_{v\epsilon}(t) = R_{12} \left[(tI_p - R_{22})^2 + \epsilon^2 I_p \right]^{-1} R_{21},$$

for $\epsilon > 0$. Observe that $Q_v > Q_{v\epsilon}(\hat{\ell}_v)$, such that

$$\liminf \lambda_{\min}(Q_v) \geq \liminf \lambda_{\min}(Q_{v\epsilon}(\hat{\ell}_v)) = \liminf \lambda_{\min}(Q_{v\epsilon}(b_\gamma) + \Delta_{v\epsilon}),$$

where $\Delta_{v\epsilon} := Q_{v\epsilon}(\hat{\ell}_v) - Q_{v\epsilon}(b_\gamma)$ (Recall that $\hat{\ell}_v \xrightarrow{\text{a.s.}} b_\gamma$). We show that $\Delta_{v\epsilon} \xrightarrow{\text{a.s.}} 0$, and

$$Q_{v\epsilon}(b_\gamma) \xrightarrow{\text{a.s.}} \int x \left[(b_\gamma - x)^2 + \epsilon^2 \right]^{-1} F_\gamma(dx) \cdot \Gamma = c_\gamma(\epsilon) \Gamma, \tag{6.38}$$

say. Because $\lambda_{\min}(\cdot)$ is a continuous function on $m \times m$ matrices, we conclude that

$$\liminf \lambda_{\min}(Q_v) \geq c_\gamma(\epsilon) \lambda_{\min}(\Gamma), \tag{6.39}$$

and because $c_\gamma(\epsilon) \rightarrow \alpha(b_\gamma + \epsilon)$ and $\alpha(b_\gamma + \epsilon) \nearrow \infty$ as $\epsilon \searrow 0$, by [JY, Appendix A], we obtain $\lambda_{\min}(Q_v) \xrightarrow{\text{a.s.}} \infty$. We write $Q_{v\epsilon}(t) = n^{-1} \bar{X}_1 \check{B}_{n\epsilon}(t) \bar{X}_1$, with

$$\begin{aligned}\check{B}_{n\epsilon}(t) &= n^{-1} \bar{X}_2^T \left[(tI_p - n^{-1} \bar{X}_2 \bar{X}_2^T)^2 + \epsilon^2 I_p \right]^{-1} \bar{X}_2 \\ &= H \text{diag}\{f_\epsilon(\mu_i, t)\} H^T,\end{aligned}$$

if we write the singular-value decomposition of $n^{-1/2} \bar{X}_2 = V \mathcal{M}^{1/2} H^T$, with $\mathcal{M} = \text{diag}(\mu_i)_{i=1}^p$ and define $f_\epsilon(\mu, t) = \mu \left[(t - \mu)^2 + \epsilon^2 \right]^{-1}$. Evidently, $\|\check{B}_{n\epsilon}(t)\| \leq \epsilon^{-2} \mu_1$ is bounded almost surely. Thus, Lemma 2 may be applied to $Q_{v\epsilon}(b_\gamma)$, and because

$$n^{-1} \text{tr} \check{B}_{n\epsilon}(b_\gamma) \xrightarrow{\text{a.s.}} \int f_\epsilon(x, b_\gamma) F_\gamma(dx) = c_\gamma(\epsilon)$$

from (5.19), our claim (6.38) follows.

Now consider v_ϵ . Fix $a \in \mathbb{R}^m$ such that $\|a\|_2 = 1$, and set $b = n^{-1/2} H^T \bar{X}_1^T a$. We have

$$a^T \Delta_{v\epsilon} a = \sum_{i=1}^p b_i^2 [f_\epsilon(\mu_i, \hat{\ell}_v) - f_\epsilon(\mu_i, b_\gamma)].$$

Because $|\partial f_\epsilon(\mu, t) / \partial t| = 2\mu(t - \mu) / \left[(t - \mu)^2 + \epsilon^2 \right]^2 \leq \mu / \epsilon^3$, for $\mu, \epsilon > 0$, by the arithmetic mean–geometric-mean inequality, we have

$$|a^T \Delta_{v\epsilon} a| \leq \mu_1 \epsilon^{-3} |\hat{\ell}_v - b_\gamma| \cdot \|b\|_2^2 = \mu_1 \epsilon^{-3} |\hat{\ell}_v - b_\gamma| a^T R_{11} a \leq \mu_1 \epsilon^{-3} |\hat{\ell}_v - b_\gamma| \hat{\ell}_1 \xrightarrow{\text{a.s.}} 0,$$

from Cauchy's interlacing inequality for eigenvalues of symmetric matrices, Theorem 1-(i) and Theorem 3-(i). Therefore, $\Delta_{v\epsilon} \xrightarrow{\text{a.s.}} 0$, and the proof of (6.39) and, hence, of Theorem 3-(ii) is complete.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgments

This work was supported, in part, by NIH R01 EB001988 (IMJ, JY), the Hong Kong RGC General Research Fund 16202918 (MRM, DMJ), and a Samsung Scholarship (JY).

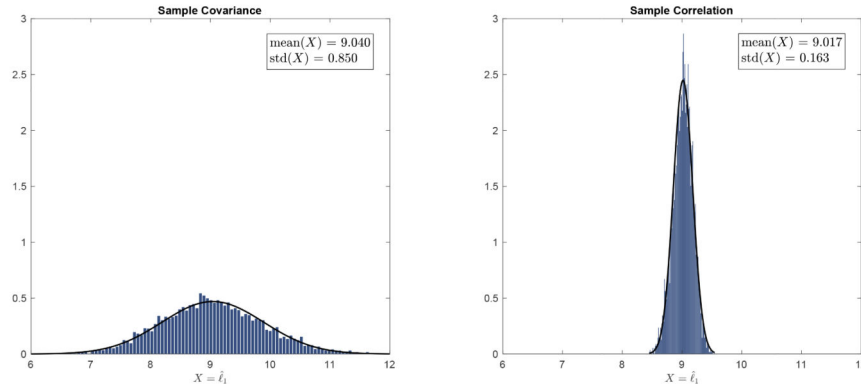
References

- Bai Z and Yao J-F (2008). Central limit theorems for eigenvalues in a spiked population model. *Annales de l'Institut Henri Poincaré, Probabilités et Statistiques* 44(3), 447–474.
- Bai ZD and Silverstein J (2009). *Spectral Analysis of Large Dimensional Random Matrices* (2nd ed.). New York: Springer.
- Baik J, Ben Arous G, and Pécché S (2005). Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices. *Annals of Probability* 33(5), 1643–1697.

- Baik J and Silverstein JW (2006). Eigenvalues of large sample covariance matrices of spiked population models. *Journal of Multivariate Analysis* 97(6), 1382–1408.
- Bao Z, Pan G, and Zhou W (2012). Tracy-Widom law for the extreme eigenvalues of sample correlation matrices. *Electronic Journal of Probability* 17, 1–32.
- Benaych-Georges F and Nadakuditi RR (2011). The eigenvalues and eigenvectors of finite, low rank perturbations of large random matrices. *Advances in Mathematics* 227(1), 494–521.
- Bianchi P, Najim J, Maida M, and Debbah M (2009). Performance analysis of some eigen-based hypothesis tests for collaborative sensing. In *2009 IEEE/SP 15th Workshop on Statistical Signal Processing*, pp. 5–8.
- Bloemendal A, Knowles A, Yau H-T, and Yin J (2016). On the principal components of sample covariance matrices. *Probability Theory and Related Fields* 164(1), 459–552.
- Boik RJ (2003). Principal component models for correlation matrices. *Biometrika* 90(3), 679–701.
- Cai TT and Jiang T (2011). Limiting laws of coherence of random matrices with applications to testing covariance structure and construction of compressed sensing matrices. *Annals of Statistics* 39(3), 1496–1525.
- Cai TT and Jiang T (2012). Phase transition in limiting distributions of coherence of high-dimensional random matrices. *Journal of Multivariate Analysis* 107, 24–39.
- Cocco S, Monasson R, and Sessak V (2011). High-dimensional inference with the generalized Hopfield model: Principal component analysis and corrections. *Physical Review E* 83(5), 051123.
- Cocco S, Monasson R, and Weigt M (2013). From principal component to direct coupling analysis of co-evolution in proteins: Low-eigenvalue modes are needed for structure prediction. *PLoS Computational Biology* 9(8), 1–17.
- Cochran D, Gish H, and Sinno D (1995). A geometric approach to multiple-channel signal detection. *IEEE Transactions on Signal Processing* 43(9), 2049–2057.
- Couillet R and Debbah M (2011). *Random Matrix Methods for Wireless Communications*. Cambridge University Press.
- Couillet R and Hachem W (2013). Fluctuations of spiked random matrix models and failure diagnosis in sensor networks. *IEEE Transactions on Information Theory* 59(1), 509–525.
- Dahirel V, Shekhar K, Pereyra F, Miura T, Artyomov M, Talsania S, Allen TM, Altfeld M, Carrington MN, Irvine DJ, Walker BD, and Chakraborty AK (2011). Coordinate linkage of HIV evolution reveals regions of immunological vulnerability. *Proceedings of the National Academy of Sciences* 108(28), 11530–11535.
- El Karoui N (2009). Concentration of measure and spectra of random matrices: Applications to correlation matrices, elliptical distributions and beyond. *Annals of Applied Probability* 19(6), 2362–2405.
- Fang C and Krishnaiah P (1982). Asymptotic distributions of functions of the eigenvalues of some random matrices for nonnormal populations. *Journal of Multivariate Analysis* 12(1), 39–63.
- Gao J, Han X, Pan G, and Yang Y (2017). High dimensional correlation matrices: The central limit theorem and its applications. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 79(3), 677–693.
- Girshick MA (1939). On the sampling theory of roots of determinantal equations. *Annals of Mathematical Statistics* 10(3), 203–224.
- Hachem W, Loubaton P, Mestre X, Najim J, and Vallet P (2013). A subspace estimator for fixed rank perturbations of large random matrices. *Journal of Multivariate Analysis* 114, 427–447.
- Hero A and Rajaratnam B (2011). Large-scale correlation screening. *Journal of the American Statistical Association* 106(496), 1540–1552.
- Hero A and Rajaratnam B (2012). Hub discovery in partial correlation graphs. *IEEE Transactions on Information Theory* 58(9), 6064–6078.
- Jiang T (2004a). The asymptotic distributions of the largest entries of sample correlation matrices. *Annals of Applied Probability* 14(2), 865–880.
- Jiang T (2004b). The limiting distributions of eigenvalues of sample correlation matrices. *Sankhy : The Indian Journal of Statistics* (2003–2007) 66(1), 35–48.

- Johnstone IM (2001). On the distribution of the largest eigenvalue in principal components analysis. *Annals of Statistics* 29(2), 295–327.
- Johnstone IM and Yang J (2018). Notes on asymptotics of sample eigenstructure for spiked models with non-Gaussian data. arXiv:1810.10427.
- Kollo T and Neudecker H (1993). Asymptotics of eigenvalues and unit-length eigenvectors of sample variance and correlation matrices. *Journal of Multivariate Analysis* 47(2), 283–300.
- Konishi S (1979). Asymptotic expansions for the distributions of statistics based on the sample correlation matrix in principal component analysis. *Hiroshima Mathematical Journal* 9(3), 647–700.
- Leshem A and van der Veen A-J (2001). Multichannel detection of Gaussian signals with uncalibrated receivers. *IEEE Signal Processing Letters* 8(4), 120–122.
- Liu H, Hu Z, Mian A, Tian H, and Zhu X (2014). A new user similarity model to improve the accuracy of collaborative filtering. *Knowledge-Based Systems* 56, 156–166.
- Mestre X and Vallet P (2017). Correlation tests and linear spectral statistics of the sample correlation matrix. *IEEE Transactions on Information Theory* 63(7), 4585–4618.
- Paul D (2007). Asymptotics of sample eigenstructure for a large dimensional spiked covariance model. *Statistica Sinica* 17, 1617–1642.
- Pillai NS and Yin J (2012). Edge universality of correlation matrices. *Annals of Statistics* 40(3), 1737–1763.
- Perou V, Gopikrishnan P, Rosenow B, Amaral L, Guhr T, and Stanley H (2002). A random matrix approach to cross-correlations in financial data. *Physical Review E* 65, 066126.
- Quadeer AA, Louie RHY, Shekhar K, Chakraborty AK, Hsing I-M, and McKay MR (2014). Statistical linkage analysis of substitutions in patient-derived sequences of genotype 1a Hepatitis C virus nonstructural protein 3 exposes targets for immunogen design. *Journal of Virology* 88(13), 7628–7644. [PubMed: 24760894]
- Quadeer AA, Morales-Jimenez D, and McKay MR (2018). Co-evolution networks of HIV/HCV are modular with direct association to structure and function. *PLoS Computational Biology* 14(9), 1–29.
- Ruan D, Meng T, and Gao K (2016). A hybrid recommendation technique optimized by dimension reduction. In 2016 8th International Conference on Modelling, Identification and Control (ICMIC), pp. 429–433.
- Schott JR (1991). A test for a specific principal component of a correlation matrix. *Journal of the American Statistical Association* 86(415), 747–751.
- Vallet P, Mestre X, and Loubaton P (2015). Performance analysis of an improved MUSIC DoA estimator. *IEEE Transactions on Signal Processing* 63(23), 6407–6422.
- Xiao H and Zhou W (2010). Almost sure limit of the smallest eigenvalue of some sample correlation matrices. *Journal of Theoretical Probability* 23(1), 1–20.
- Yang L, McKay MR, and Couillet R (2018). High-dimensional MVDR beamforming: Optimized solutions based on spiked random matrix models. *IEEE Transactions on Signal Processing* 66(7), 1933–1947.
- Yao J, Zheng S, and Bai Z (2015). *Large Sample Covariance Matrices and High-Dimensional Data Analysis*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press.

(a) Histogram of the largest sample eigenvalue



(b) Scatter plot of sample-to-population eigenvector projections

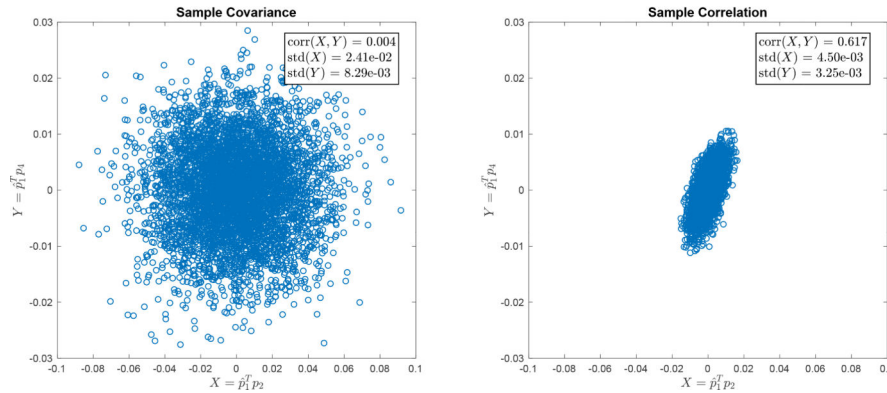
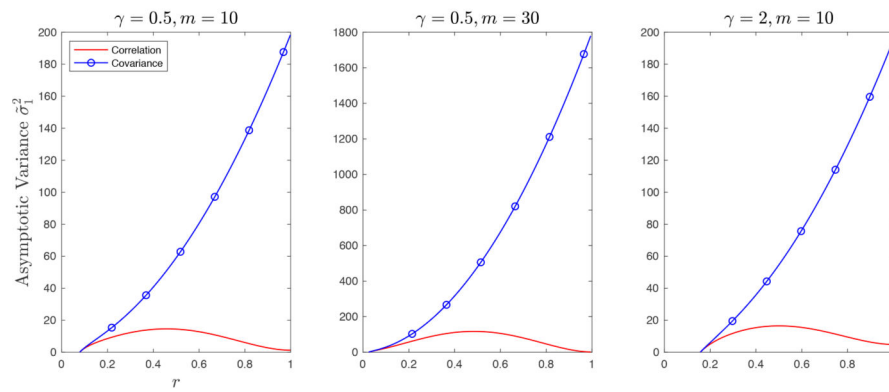


Figure 1:

A simple simulation shows remarkable distributional differences between sample covariance and sample correlation. From $n = 200$ i.i.d. Gaussian samples, $x_i \in \mathbb{R}^{100}$, with covariance $\Sigma = \text{blkdiag}(\Sigma_s, I_{90})$, where $(\Sigma_s)_{i,j}^{10} = (r^{|i-j|})_{i,j=1}^{10}$, for $r = 0.95$, we compute the sample covariance and sample correlation, and show: (a) the empirical density (normalized histogram) of the largest sample eigenvalue, along with a Gaussian distribution with its estimated mean and standard deviation (solid line), and (b) a scatter plot of the leading sample eigenvector, projected onto the second (x-axis) and fourth (y-axis) population eigenvectors. A striking variance reduction is observed in the sample correlation for both (a) and (b). A similar variance reduction is observed for different choices of population eigenvectors in (b); the selected choice (being the second and fourth eigenvectors) facilitates the illustration of an additional correlation effect in the sample-to-population eigenvector projections.

(a) Largest sample eigenvalue $\hat{\ell}_1$



(b) Sample-to-population eigenvector projection $p_2^T a_1$

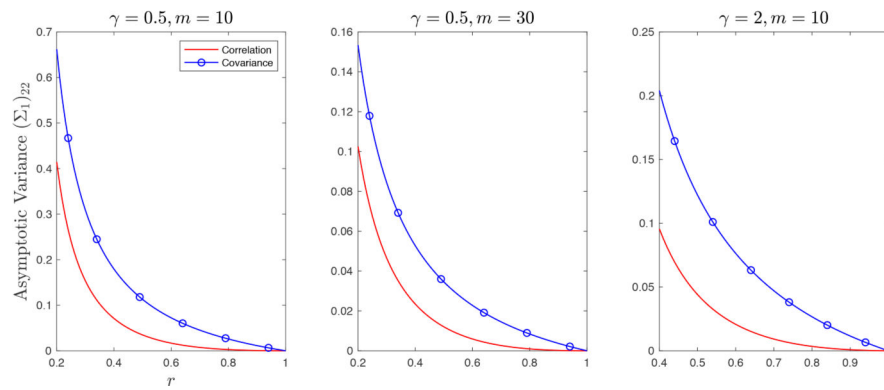


Figure 2: Differences in the fluctuations of sample eigenvalues and eigenvectors for an example Gaussian model with $\Gamma = (1 - r)I_m + r1_m1_m^T$. Asymptotic variances are shown for (a) the largest sample eigenvalue $\hat{\ell}_1$, and (b) the normalized sample-to-population eigenvector projection $p_2^T a_1$.