# Convergence of Sample Eigenvalues, Eigenvectors, and Principal Component Scores for Ultra-High Dimensional Data

**Seunggeun Lee**,
Department of Biostatistics, University of Michigan, 1415 Washington Heights, Ann Arbor, Michigan 48109, U.S.A

**Fei Zou**, and
Department of Biostatistics, University of North Carolina, 135 Dauer Drive, Chapel Hill, North Carolina 27599, U.S.A

**Fred A. Wright**
Bioinformatics Research Center, North Carolina State University, 1 Lampe Drive, Raleigh, North Carolina 27607, U.S.A

Seunggeun Lee: leeshawn@umich.edu; Fei Zou: fzou@bios.unc.edu; Fred A. Wright: fred_wright@ncsu.edu

## Summary

The development of high-throughput biomedical technologies has led to increased interest in the analysis of high-dimensional data where the number of features is much larger than the sample size. In this paper, we investigate principal component analysis under the ultra-high dimensional regime, where both the number of features and the sample size increase as the ratio of the two quantities also increases. We bridge the existing results from the finite and the high-dimension low sample size regimes, embedding the two regimes in a more general framework. We also numerically demonstrate the universal application of the results from the finite regime.

### Some key words

High-Dimension Low Sample Size Data; Principal Component Analysis; Random Matrix

## 1. Introduction

With the development of modern high-throughput technologies, it is common to encounter data with many more features, $p$, than the number of samples, $n$. In modern genomics applications, for instance, the number of features often ranges from tens of thousands to millions, while the corresponding sample sizes typically range from hundreds to thousands. For those high-dimensional data, principal component analysis is popular for data exploration and dimension reduction. Since principal component analysis is based on the eigenvalues and eigenvectors of the sample covariance matrix, its performance largely depends on the behavior of the sample eigenvalues and eigenvectors.

In their seminal paper on random matrices, Marčenko & Pastur (1967) derived the asymptotic distribution of the sample eigenvalues under the finite $\gamma$ regime, where $p \to \infty$, $n \to \infty$ and $p/n \to \gamma < \infty$. Specifically, they showed that the sample eigenvalues follow the Marčenko–Pastur law when all the population eigenvalues are identical. For data where the true signal is embedded in a low dimensional space, Johnstone (2001) introduced the spiked eigenvalue model, where a small number of population eigenvalues are substantially larger than the rest. Under this model, asymptotic results on the sample eigenvalues and eigenvectors have been derived (Baik & Silverstein, 2006; Paul, 2007; Nadler, 2008; Lee et al., 2010) for the finite $\gamma$ asymptotic regime.

These results are useful for evaluating the performances of principal component analysis (Lee et al., 2010). However, one may be concerned about the applicability of the theoretical results from the finite $\gamma$ regime to ultra-high dimensional data, such as next generation sequencing data, where millions of genetic variants are collected from tens or a few hundreds of samples. Addressing this question is urgent, as the availability of such ultra-high dimensional genomic datasets is expected to increase as the cost of high-throughput technologies decreases. In this paper, we derive asymptotic results that provide theoretical justification for applying the results from the finite $\gamma$ regime to ultra-high dimensional data. In addition, we compare our results to those from the high-dimension low sample size regime (Hall et al., 2005; Ahn et al., 2007; Jung & Marron, 2009; Jung et al., 2012).

The finite $\gamma$ and the high-dimension low sample size regimes are based on two seemingly disparate assumptions. In the high-dimension low sample size regime, $n$ is treated as fixed and the population eigenvalues increase with rate $p^a$. In the finite $\gamma$ regime, the population eigenvalues are assumed to be fixed but $n$ grows with $p$ at a constant rate. Our new results on the ultra-high dimensional regime bridge the asymptotic results from the two extreme regimes and improve our understanding of principal component analysis on high-dimensional data.

## 2. Method

### 2.1. General Setting

Throughout this paper, we assume that $n$ is a function of $p$, and denote it by $n_p$ whenever needed. We further define $\gamma_p = p/n_p$. Let $\sum_p = E_p \Lambda_p E_p^T$ be a $p \times p$ nonnegative matrix with an ordered eigenvalue matrix $\Lambda_p = \mathrm{diag}(\lambda_{p1}, \ldots, \lambda_{pp})$ and an orthogonal eigenvector matrix $E_p = (e_{p1}, \ldots, e_{pp})$. Both eigenvalues and eigenvectors are fixed sequences which depend on $p$. Define the $p \times n$ data matrix, $X_p = E_p \Lambda_p^{1/2} Z_p$, where $Z_p$ is a $p \times n$ random matrix whose elements $z_{ij}$ are independent and identically distributed with $E(Z_{ij})=0$, $E(Z_{ij}^2)=\sigma^2$ and $E(z_{ij}^4)<\infty$. The sample covariance matrix $S_p$ equals

$$S_p = n^{-1} X_p X_p^T = n^{-1} E_p \Lambda_p^{1/2} Z_p Z_p^T \Lambda_p^{1/2} E_p^T,$$

and the corresponding population covariance matrix of $X_p$ is $\sigma^2 \Sigma_p$. The $\sigma^2 \lambda_{pv}$s are the underlying population eigenvalues. The spectral decomposition of the sample covariance matrix is $S_p = U_p D_p U_p^T$, where $D_p = \mathrm{diag}(d_{p1}, \ldots, d_{pp})$ is the diagonal matrix of the ordered sample eigenvalues, and $U_p = (u_{p1}, \ldots, u_{pp})$ is the corresponding $p \times p$ sample eigenvector matrix. The $v$th sample principal component score vector, $\widehat{p_v} = (\widehat{p_{v1}}, \ldots, \widehat{p_{vn}})^T$, equals $X^T u_v$. For a new sample with variable $x_{\text{new}}$, its $v$th predicted principal component score is $\hat{q}_v = X_{\text{new}}^T u_v$. Before introducing the main results, we define additional notation for the remainder of the paper. Suppose $a_p$ and $b_p$ are two sequences. We write $a_p \asymp b_p$ if $a_p = O(b_p)$ and $b_p = O(a_p)$, and $a_p \ll b_p$ if $a_p/b_p = o(1)$. For simplicity, we hence suppress the subscript $p$ unless we wish to emphasize a quantity's dependence on $p$, except for the population eigenvector matrix, which is always denoted by $E_p$.

## 2.2. Main Results

In the sequel, we assume $\gamma_p \to \infty$ and $n \to \infty$ as $p \to \infty$. We further assume the spiked eigenvalue model (Johnstone, 2001) in which the first $m$ population eigenvalues are substantially larger than the remaining non-spiked eigenvalues. In the random matrix context, it is typically assumed that all non-spiked population eigenvalues equal unity (Johnstone, 2001; Baik & Silverstein, 2006). This strong condition is unlikely to be satisfied in many situations. We define two weaker sphericity conditions. Let

$\phi(k) = (p - m)^{-1} \sum_{v=m+1}^{p} (\lambda_v - \overline{\lambda})^k$ be the $k$th central moment of the non-spiked

population eigenvalues, where $\overline{\lambda} = (p - m)^{-1} \sum_{v=m+1}^{p} \lambda_v$.

*Condition* 1. The non-spiked population eigenvalues satisfy $\varphi(2) = o(n^{-2}p)$.

*Condition* 2. The non-spiked population eigenvalues satisfy $\varphi(2) = o(n^{-3/2}p)$, $\varphi(4) = O(1)$, and $\varphi(4) = o(n^{-4}p^3)$.

Condition 1 is closely related to the sphericity measure in John (1971, 1972) and the ($\epsilon_m$ condition of Jung & Marron (2009). Detailed explanations of both conditions can be found in the Supplementary Material. The following theorem summarizes the convergence results of the sample eigenvalues and eigenvectors.

**Theorem 1**—Let $c_v = \lambda_v/\gamma_p$ ($v \quad m$). Suppose that $c_m < \cdots < c_1$, and $c_m \asymp \cdots \asymp c_1$. Let the remaining population eigenvalues satisfy Condition 1 or Condition 2.

**i.** When $c_v$ is bounded away from zero, for $v \quad m$, $d_v \lambda_v^{-1} - \sigma^2 c_v^{-1}(c_v+1) \to 0$ in probability, and $|\langle e_v, u_v \rangle| - \{c_v(c_v + 1)^{-1}\}^{1/2} \to 0$ in probability, where $\langle . \rangle$ is the inner product between two vectors. For $v > m$, $d_v \gamma_p^{-1} \to \sigma^2$ in probability.

**ii.** ii) When $c_v = o(1)$, for all $v$, $d_v \gamma_p^{-1} \to \sigma^2$ in probability, and $|\langle e_v, u_v \rangle| \to 0$ in probability.

The proof can be found in the Supplementary Material. Theorem 1 includes convergence results for both the spiked and non-spiked sample eigenvalues. These results clearly indicate that the asymptotic behavior of sample eigenvalues and eigenvectors depends on $c_v$, which

can be viewed as a signal to noise ratio, where $\lambda_v$ represents the signal strength and $\gamma_p$ serves as a surrogate of the noise level.

When $\lambda_v$ grows at the same rate as, or at a higher rate than, $\gamma_p$, the spiked eigenvalues are separable from the bulk. When $\lambda_v$ grows at a slower rate than $\gamma_p$, i.e., $c_v = o(1)$, the spiked eigenvalues cannot be separated from the non-spiked eigenvalues. Theorem 1 also shows that $d_v$ is inconsistent. The sample eigenvectors show a similar pattern. Examples on the asymptotic behavior of the sample eigenvalues and eigenvectors under several conditions are described in the Supplementary Material. To mimic the high-dimension low sample size regime, let $\lambda_v$ be a function of $\gamma_p^\alpha$ such that a limit of $\lambda_v/\gamma_p^\alpha$ exists and is finite. Now we have the following corollary.

**Corollary 1**—*Let* $\lambda_v/\gamma_p^\alpha \to \tilde{c}_v$ $(v \leq m), \tilde{c}_m < \cdots < \tilde{c}_1$, *and the remaining population eigenvalues satisfy Condition 1 or Condition 2. Then, for $v \quad m$, $d_v/\max(\gamma_p, \gamma_p^\alpha)$ converges in probability to $\sigma^2\tilde{c}_v$, $\sigma^2\tilde{c}_v + \sigma^2$, and $\sigma^2$ when $\alpha > 1$, $\alpha = 1$, and $\alpha < 1$, respectively. With the same assumption, $|\langle e_v, u_v \rangle|$ converges in probability to unity, $\{\tilde{c}_v(\tilde{c}_v + 1)^{-1}\}^{1/2}$, and zero when $\alpha > 1$, $\alpha = 1$, and $\alpha < 1$, respectively.*

The proof can be found in the Supplementary Material. The corollary allows us to compare our results to those from the high-dimension low sample size regime. See Section 2.3 for details. After principal component analysis, the sample principal component scores are often used to summarize data. Predicted principal component scores may also be calculated on new samples for a variety of reasons (Jolliffe, 2002). The next theorem presents the asymptotic results on the principal component scores under the ultra-high dimensional regime.

**Theorem 2**—Suppose the assumptions in Theorem 1 hold, and $c_v$ $(v \quad m)$ is bounded away from zero. Let $p_v = X^T e_v$ be the $v$th population principal component score derived from the corresponding $v$th population eigenvector, and corr$(\cdot, \cdot)$ be the correlation function. Then, for $v \quad m$, corr$(p_v, \hat{p_v}) \to 1$ in probability, and $\left\{ E(\hat{q}_v^2)/E(\hat{p}_{vj}^2) \right\}^{1/2} - c_v(c_v+1)^{-1} \to 0$ inprobability, for all $j = 1, \ldots, n$.

The proof is given in the Supplementary Material. One striking feature in Theorem 2 is that the correlation between $p_v$ and $\hat{p_v}$ can converge to unity even when the corresponding sample eigenvector is not consistent. Combining Theorems 1 and 2, we conclude that $\hat{p_v}$ can accurately estimate $p_v$ whenever its corresponding sample eigenvalue is separable from the bulk. This interesting result may partially explain the success of principal component analysis for high dimensional datasets, such as genome-wide association data (Price et al., 2006; Patterson et al., 2006). This theorem also illustrates the shrinkage phenomena of the predicted principal component scores, previously reported by Lee et al. (2010). To apply the asymptotic results in Theorem 1 and 2 to data, we need to estimate $\sigma^2$. Lee et al. (2010) proposed an algorithm to rescale the data to ensure $\sigma^2$ of the rescaled data equal to unity. The same approach can be applied to ultra-high dimensional data.

### 2.3. Comparisons to existing asymptotic results

In the finite $\gamma$ framework, it is typically assumed that the spiked population eigenvalues are finite. Under the spiked population model, the following results have been established (Baik & Silverstein, 2006; Paul, 2007; Lee et al., 2010). For sample eigenvalues,

$$
\begin{aligned}
d_v - \sigma^2 \lambda_v \left\{ 1 + \gamma(\lambda_v - 1)^{-1} \right\} &= o_p(1), \quad \lambda_v > 1 + \gamma^{1/2}, \\
d_v - \sigma^2 (1 + \gamma^{1/2})^2 &= o_p(1), \quad\quad \lambda_v \leq 1 + \gamma^{1/2},
\end{aligned}
\tag{1}
$$

and for predicted principal component scores,

$$
\left\{ E(\hat{q}_v^2)/E(\hat{p}_{vj}^2) \right\}^{1/2} - (\lambda_v - 1)(\lambda_v + \gamma - 1)^{-1} = o_p(1). \tag{2}
$$

Equation (1) shows that if $\lambda_v > 1 + \gamma^{1/2}$, its corresponding sample eigenvalue is separable from the bulk. Interestingly, the result in (1) for $\lambda_v > 1 + \gamma^{1/2}$ is equivalent to the asymptotic result in Theorem 1 when $\lambda_v$ is relatively large. To see this, note that by replacing $\lambda_v$ in (1) with $c_v \gamma$, we obtain

$$
d_v \lambda_v^{-1} \approx \sigma^2 + (\sigma^2 \gamma)(\lambda_v - 1)^{-1} \approx \sigma^2 c_v^{-1}(c_v + 1),
$$

which accords with Theorem 1. For the predicted principal component scores, the same holds true. By (2),

$$
\left\{ E(\hat{q}_v^2)/E(\hat{p}_{vj}^2) \right\}^{1/2} \approx (\lambda_v - 1)(\lambda_v + \gamma - 1)^{-1} \approx c_v(c_v + 1)^{-1},
$$

which is consistent with Theorem 2. Similar conclusions can be drawn for the sample eigenvectors and principal component scores and are omitted here. In summary, for ultra-high dimensional data, both the finite $\gamma$ and ultra-high dimensional asymptotic results can be used to investigate the behavior of sample eigenvalues, eigenvectors and principal component scores, and both produce similar conclusions.

Under the high-dimension low sample size regime, the spiked population eigenvalues are assumed to grow at rate $p^\alpha$ ($a > 0$). Jung et al. (2012) proved the following results with an additional Gaussian assumption on $X$. For the first sample eigenvalue,

$$
\frac{d_1}{\max(p, p^\alpha)} \rightarrow
\begin{cases}
\sigma^2 \hat{c}_1 \chi_n^2 / n, & \alpha > 1, \\
\sigma^2 \hat{c}_1 \chi_n^2 / n + \sigma^2 / n, & \alpha = 1, \\
\sigma^2 / n, & \alpha < 1,
\end{cases}
\tag{3}
$$

in distribution, where $\hat{c}_1 = \lim_{p \to \infty} \lambda_1 / p^\alpha$, and $\chi_n^2$ denotes the chi-square distribution with $n$ degrees of freedom. For the first sample eigenvector,

$$|\langle e_1, u_1 \rangle| \rightarrow \begin{cases} 1, & \alpha > 1, \\ \{\hat{c}_1 \chi_n^2 / (\hat{c}_1 \chi_n^2 + 1)\}^{1/2}, & \alpha = 1, \\ 0, & \alpha < 1, \end{cases} \quad (4)$$

in distribution. Results with more relaxed assumptions can be found in Jung et al. (2012). In the high-dimension low sample size regime, the asymptotic behavior of sample eigenvalues and eigenvectors depends on the relative growth rate of $\lambda_v$ over $p$. In Corollary 1, $\lambda_v$ is expressed as a function of $\gamma_p$, instead of $p$ directly. However, it should be noted that $\gamma_p^{\alpha}$ is equivalent to $p^a$ in the high-dimension low sample size regime where $n$ is treated as fixed. Equation (3) shows that when $a > 1$, the distribution of the scaled sample eigenvalue converges in distribution to the random variable $\sigma^2 \hat{c}_1 \chi_n^2 / n$ as $p \rightarrow \infty$. Combining this with the fact that $\chi_n^2 / n \rightarrow 1$ in probability as $n \rightarrow \infty$, we end up with the same conclusion in Corollary 1. When $a = 1$, $d_1 - \sigma^2 \lambda_1 \asymp \sigma^2 p/n$, and thus the sample eigenvalue is biased with a bias $\sigma^2 p/n$. Equation (4) indicates that the first sample eigenvector is consistent when $a > 1$ and is asymptotically perpendicular to the first population eigenvector when $a < 1$. When $a = 1$, the sample eigenvector is neither consistent nor asymptotically perpendicular to the first population eigenvector. In conclusion, our asymptotic results clearly parallel the asymptotic results for the high-dimension low sample size regime, and embed them within a larger framework.

## 3. Numerical Study

We conducted simulations to illustrate our theoretical results. A $p \times n$ data matrix $X$ was generated from $N(0, \Lambda)$ with $\Lambda = \text{diag}(\lambda_1, \ldots, \lambda_p)$. We set $\lambda_1 = c_1 \gamma_p$, $\lambda_2 = c_2 \gamma_p$ and $\lambda_3 = \ldots = \lambda_p = 1$. The first and second population eigenvectors were $e_1 = (1, 0, \ldots, 0)$ and $e_2 = (0, 1, 0, \ldots, 0)$, respectively. Four different sets of $c_v$s were selected to represent different scenarios: no spiked eigenvalues, $c_1 = c_2 = \gamma_p^{-1}$; very small spiked eigenvalues, $c_1 = \gamma_p^{-1/2}, c_2 = 0.7\gamma_p^{-1/2}$; moderate spiked eigenvalues, $c_1 = 1$, $c_2 = 0.7$; and very large spiked eigenvalues, $c_1 = \gamma_p^{1/2}, c_2 = 0.7\gamma_p^{1/2}$. The first two scenarios correspond to the case that $c_v = o(1)$, and the last two to the case that $c_v$ is bounded away from zero. Two different $\gamma_p$ values, 500 and 2000, were considered, and the sample size was fixed at 100. For each of the simulation setups, we generated 500 datasets and computed the sample eigenvalues and the inner products between the sample and population eigenvectors.

Table 1 reports the medians and inter-quartile ranges of the estimates. The theoretical asymptotic values of the sample eigenvalues and the inner products from the finite $\gamma$ and the ultra-high dimensional regimes are also presented. The sample eigenvalues were rescaled by $\gamma_p$. For data with no spiked or with very small spiked eigenvalues, the first and second sample eigenvalues are slightly upward-biased from unity. However, they match well with the theoretical ones from the finite $\gamma$ regime. For data with moderate or large spiked eigenvalues, the theoretical estimates from the finite $\gamma$ regime and the ultra-high dimensional regime are identical, and are well matched with the sample eigenvalues. For the inner products of the sample eigenvectors, the empirical estimates match well with the ones from

the finite $\gamma$ and ultra-high dimensional regimes, and the two sets of the theoretical results are identical.

Table 2 summarizes the results for the sample and predicted principal component scores. For the sample principal component scores, the median and inter-quartile ranges of their Pearson correlations with the population principal component scores were calculated. The theoretical results from both the finite $\gamma$ and ultra-high dimensional regimes are identical and both match well with the empirical estimates. For the predicted principal component scores, we followed exactly the same simulation procedure as described above to generate a new dataset for each of the simulated dataset. We computed the predicted principal component scores on each new dataset. The empirical shrinkage factor was calculated as the ratio of the means of the squared predicted and sample principal component scores. Again, similar conclusions hold. The theoretical results from both the finite $\gamma$ and ultra-high dimensional regimes are effectively identical. The empirical estimates approximate the theoretical results very well.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgments

## References

Ahn J, Marron JS, Muller KM, Chi YY. The high-dimension, low-sample-size geometric representation holds under mild conditions. Biometrika. 2007; 94:760–766.

Baik J, Silverstein JW. Eigenvalues of large sample covariance matrices of spiked population models. J Multivariate Anal. 2006; 97:1382–1408.

Hall P, Marron JS, Neeman A. Geometric representation of high dimension, low sample size data. J R Statist Soc B. 2005; 67:427–444.

John S. Some optimal multivariate tests. Biometrika. 1971; 58:123–127.

John S. The distribution of a statistic used for testing sphericity of normal distributions. Biometrika. 1972; 59:169–173.

Johnstone IM. On the distribution of the largest eigenvalue in principal components analysis. Ann Statist. 2001; 29:295–327.

Jolliffe, I. Principal Component Analysis. Springer; New York: 2002.

Jung S, Marron JS. PCA consistency in high dimension, low sample size context. Ann Statist. 2009; 37:4104–4130.

Jung S, Sen A, Marron JS. Boundary behavior in high dimension, low sample size asymptotics of PCA. J Multivariate Anal. 2012; 109:190–203.

Lee S, Zou F, Wright F. Convergence and prediction of principal component scores in high-dimensional settings. Ann Statist. 2010; 38:3605–3629.

Mar enko V, Pastur L. Distribution of eigenvalues for some sets of random matrices. Sbornik: Mathematics. 1967; 1:457–483.

Nadler B. Finite sample approximation results for principal component analysis: a matrix perturbation approach. Ann Statist. 2008; 36:2791–2817.

Patterson N, Price A, Reich D. Population structure and eigenanalysis. PLoS Genet. 2006; 2:e190. [PubMed: 17194218]

Paul D. Asymptotics of sample eigenstructure for a large dimensional spiked covariance model. Statist Sinica. 2007; 17:1617–1642.

Price A, Patterson N, Plenge R, Weinblatt M, Shadick N, Reich D. Principal components analysis corrects for stratification in genome-wide association studies. Nat Genet. 2006; 38:904–909. [PubMed: 16862161]

**Table 1**

Rescaled sample eigenvalues and eigenvectors based on 500 simulations. Ultra-High and Finite $\gamma$ present theoretical asymptotic values from the ultra-high dimensional and the finite $\gamma$ regimes, respectively. Observed presents the median of the estimates, with the interquartile range in parentheses.

| | Principal component | $n_p$ | Type | No Spike | Very small Spike | Moderate Spike | Very Large Spike |
|---|---|---|---|---|---|---|---|
| Eigenvalues | 1 | 500 | Ultra-High | 1.00 | 1.00 | 2.00 | 23.36 |
| | | | Finite $\gamma$ | 1.09 | 1.09 | 2.00 | 23.36 |
| | | | Observed | 1.09(0.004) | 1.09(0.004) | 2.02(0.183) | 23.6(3.868) |
| | | 2000 | Ultra-High | 1.00 | 1.00 | 2.00 | 45.72 |
| | | | Finite $\gamma$ | 1.05 | 1.05 | 2.00 | 45.72 |
| | | | Observed | 1.04(0.002) | 1.05(0.002) | 2.02(0.207) | 46.75(7.605) |
| | 2 | 500 | Ultra-High | 1.00 | 1.00 | 1.70 | 16.65 |
| | | | Finite $\gamma$ | 1.09 | 1.09 | 1.70 | 16.65 |
| | | | Observed | 1.08(0.003) | 1.09(0.003) | 1.67(0.125) | 15.98(2.680) |
| | | 2000 | Ultra-High | 1.00 | 1.00 | 1.70 | 32.31 |
| | | | Finite $\gamma$ | 1.05 | 1.05 | 1.70 | 32.31 |
| | | | Observed | 1.04(0.001) | 1.04(0.002) | 1.68(0.131) | 30.96(5.370) |
| Eigenvectors | 1 | 500 | Ultra-High | 0.00 | 0.00 | 0.71 | 0.98 |
| | | | Finite $\gamma$ | 0.00 | 0.00 | 0.71 | 0.98 |
| | | | Observed | 0.00(0.004) | 0.07(0.069) | 0.69(0.056) | 0.96(0.061) |
| | | 2000 | Ultra-High | 0.00 | 0.00 | 0.71 | 0.99 |
| | | | Finite $\gamma$ | 0.00 | 0.00 | 0.71 | 0.99 |
| | | | Observed | 0.00(0.002) | 0.05(0.048) | 0.69(0.056) | 0.97(0.054) |
| | 2 | 500 | Ultra-High | 0.00 | 0.00 | 0.64 | 0.97 |
| | | | Finite $\gamma$ | 0.00 | 0.00 | 0.64 | 0.97 |
| | | | Observed | 0.00(0.003) | 0.02(0.030) | 0.607(0.055) | 0.95(0.059) |
| | | 2000 | Ultra-High | 0.00 | 0.00 | 0.64 | 0.98 |
| | | | Finite $\gamma$ | 0.00 | 0.00 | 0.64 | 0.98 |
| | | | Observed | 0.00(0.002) | 0.02(0.022) | 0.61(0.052) | 0.96(0.053) |

**Table 2**

Sample and predicted principal component scores based on 500 simulations. UltraHigh and Finite $\gamma$ present theoretical asymptotic values from the ultra-high dimensional and the finite $\gamma$ regimes, respectively. Observed presents the median of the estimates, with the interquartile range in parentheses.

| | $\gamma_p$ | Type | Principal Component 1 | | Principal Component 2 | |
| --- | --- | --- | --- | --- | --- | --- |
| | | | Moderate Spike | Very Large Spike | Moderate Spike | Very Large Spike |
| Principal Component Scores | 500 | Ultra-High | 1.00 | 1.00 | 1.00 | 1.00 |
| | | Finite $\gamma$ | 1.00 | 1.00 | 1.00 | 1.00 |
| | | Observed | 0.99(0.042) | 0.99(0.043) | 0.97(0.090) | 0.97(0.087) |
| | 2000 | Ultra-High | 1.00 | 1.00 | 1.00 | 1.00 |
| | | Finite $\gamma$ | 1.00 | 1.00 | 1.00 | 1.00 |
| | | Observed | 0.99(0.039) | 0.99(0.038) | 0.97(0.078) | 0.97(0.079) |
| Shrinkage Factors | 500 | Ultra-High | 0.50 | 0.96 | 0.41 | 0.94 |
| | | Finite $\gamma$ | 0.50 | 0.96 | 0.41 | 0.94 |
| | | Observed | 0.49(0.045) | 0.93(0.120) | 0.42(0.046) | 0.97(0.122) |
| | 2000 | Ultra-High | 0.50 | 0.98 | 0.41 | 0.97 |
| | | Finite $\gamma$ | 0.50 | 0.98 | 0.41 | 0.97 |
| | | Observed | 0.49(0.045) | 0.95(0.127) | 0.42(0.041) | 1.00(0.118) |