

**Web-based Supplementary Materials for “Asymptotic conditional singular
value decomposition for high-dimensional genomic data”**

Jeffrey T. Leek

Johns Hopkins Bloomberg School of Public Health

Baltimore, MD 21205-2179

jleek@jhsph.edu

Web Appendix A

Proof of Theorem 1

\mathbf{W}^m can be broken into five terms:

$$\begin{aligned} \mathbf{W}^m &= \frac{1}{m} \mathbf{X}^{mT} \mathbf{X}^m - \hat{\sigma}_{ave}^2 \mathbf{I} \\ &= \underbrace{\frac{1}{m} \mathbf{G}^T \mathbf{\Gamma}^{mT} \mathbf{\Gamma}^m \mathbf{G}}_i + \underbrace{\frac{1}{m} \mathbf{G}^T \mathbf{\Gamma}^{mT} \mathbf{U}^m}_{ii} + \underbrace{\frac{1}{m} \mathbf{U}^{mT} \mathbf{\Gamma}^m \mathbf{G}}_{iii} + \underbrace{\frac{1}{m} \mathbf{U}^{mT} \mathbf{U}^m}_{iv} - \underbrace{\hat{\sigma}_{ave}^2 \mathbf{I}}_v \quad (\text{A.1}) \end{aligned}$$

We consider each of these terms individually.

i: This term converges to $\mathbf{G}^T \mathbf{\Delta} \mathbf{G}$ by assumption 3.

ii: Let $\mathbf{M} = \frac{1}{m} \mathbf{G}^T \mathbf{\Gamma}^{mT} \mathbf{U}^m = \frac{1}{m} \mathbf{B}^m \mathbf{U}^m$, then $m_{ij} = \frac{1}{m} \sum_{\ell=1}^m b_{i\ell} u_{\ell j}$ where $E(b_{i\ell} u_{\ell j}) = 0$ and $\text{var}(b_{i\ell} u_{\ell j}) = b_{i\ell}^2 \sigma_\ell^2$. So by the Kolmogorov Strong Law of Large Numbers (KSLLN) (Feller, 1968) $m_{ij} \rightarrow_{a.s.} 0$ for all i, j .

iii: By symmetry, this term also converges almost surely to zero.

iv: Let $\mathbf{S} = \frac{1}{m} \mathbf{U}^{mT} \mathbf{U}^m$, and consider the off-diagonal element $s_{ij} = \frac{1}{m} \sum_{\ell=1}^m u_{\ell i} u_{\ell j}$, where $E(u_{\ell i} u_{\ell j}) = 0$ and $\text{var}(u_{\ell i} u_{\ell j}) = E(u_{\ell i}^2 u_{\ell j}^2) - E(u_{\ell i} u_{\ell j})^2 = (\sigma_\ell^2)^2$. So again by KSLLN $s_{ij} \rightarrow_{a.s.} 0$. Now consider the diagonal elements $s_{ii} = \frac{1}{m} \sum_{\ell=1}^m u_{\ell i}^2$, where $E(u_{\ell i}^2) = \sigma_\ell^2$ and $\text{var}(u_{\ell i}^2) = E(u_{\ell i}^4) - E(u_{\ell i}^2)^2$. By assumptions 1 the variances are bounded, so by KSLLN $s_{ii} \rightarrow_{a.s.} \bar{\sigma}^2 = \lim_{m \rightarrow \infty} \frac{1}{m} \sum_{i=1}^m \sigma_i^2$, which exists because σ_i^2 is bounded for all i .

Combining terms (i-iv) and applying Slutsky's theorem yields: $\frac{1}{m} \mathbf{X}^{mT} \mathbf{X}^m \rightarrow_{a.s.} \mathbf{G}^T \mathbf{\Delta} \mathbf{G} +$

$\bar{\sigma}^2 \mathbf{I}$. Since the eigenvalues of a matrix are defined as roots of a determinant depending on the elements of that matrix, and since the roots of a polynomial equation are a continuous multi-valued function of the coefficients (Henriksen and Isbell, 1953), the eigenvalue function is continuous. The eigenvalues of $\frac{1}{m} \mathbf{X}^{mT} \mathbf{X}^m$ converge almost surely to the eigenvalues of $\mathbf{G}^T \Delta \mathbf{G} + \bar{\sigma}^2 \mathbf{I}$ from the continuous mapping theorem. The eigenvalues of $\mathbf{G}^T \Delta \mathbf{G} + \bar{\sigma}^2 \mathbf{I}$ are equal to $\lambda_1 + \bar{\sigma}^2, \dots, \lambda_n + \bar{\sigma}^2$; but $\lambda_{r+1}, \dots, \lambda_n$ are equal to zero by assumption, so the last $n - r$ eigenvalues consistently estimate $\bar{\sigma}^2$.

$v : \hat{\sigma}_{ave}^2 = \frac{1}{m} \sum_{i=1}^m \frac{1}{(n-\kappa)} \sum_{j=1}^n (x_{ij} - \sum_{k=1}^{\kappa} \hat{\gamma}_{ik} \hat{u}_{kj})^2 = \frac{1}{n-\kappa} \sum_{k=\kappa}^n \lambda_k(\mathbf{Z}^m)$, where $\lambda_k(\mathbf{Z}^m)$ is the k th eigenvalue of $\mathbf{Z}_m = \frac{1}{m} \mathbf{X}^{mT} \mathbf{X}^m$. But for all $\kappa > r$, $\lambda_k(\mathbf{Z}^m)$ converges to $\bar{\sigma}^2$ almost surely.

Combining terms (i-v) and applying Slutsky's theorem yields: $\mathbf{W}^m \rightarrow_{a.s.} \mathbf{G}^T \Delta \mathbf{G}$. By the same argument as above, the eigenvalues of \mathbf{W}^m converge almost surely to the eigenvalues of $\mathbf{G}^T \Delta \mathbf{G}$ from the continuous mapping theorem. Further, since both the matrix \mathbf{W}^m and the eigenvalues converge almost surely, and the eigenvectors can be obtained from a linear operation of these two elements, the eigenvectors of \mathbf{W}^m corresponding to the unique eigenvalues must converge to the corresponding eigenvectors of $\mathbf{G}^T \Delta \mathbf{G}$.

Proof of Lemma 1

We wish to show that the indicator:

$$1 \{ \lambda_k(\mathbf{W}^m) \geq c_m \} = 1 \left\{ \frac{1}{c_m} \lambda_k(\mathbf{W}^m) \geq 1 \right\}.$$

consistently distinguishes between zero and non-zero eigenvalues. If $\lambda_k(\mathbf{W}^m) = \lambda_k + O_P(m^{-\frac{1}{2}})$ then for any $c_m = O(m^{-\eta})$, $0 < \eta < \frac{1}{2}$, when $\lambda_k = 0$, $\left(\frac{1}{c_m} \lambda_k(\mathbf{W}^m)\right) = O(m^\eta) O_P(m^{-\frac{1}{2}}) = m^{\eta-\frac{1}{2}} O_P(1)$ and $1 \left\{ \frac{1}{c_m} \lambda_k(\mathbf{W}^m) \geq 1 \right\} \rightarrow_P 0$. When $\lambda_k > 0$, $\left(\frac{1}{c_m} \lambda_k(\mathbf{W}^m)\right) = O(m^\eta) \left\{ \lambda_k + O_P(m^{-\frac{1}{2}}) \right\} = O(m^\eta) + o_P \left\{ m^{\eta-\frac{1}{2}} \right\} \rightarrow \infty$. So when $\lambda_k > 0$ and $0 < \eta < \frac{1}{2}$, $1 \left\{ \frac{1}{c_m} \lambda_k(\mathbf{W}^m) \geq 1 \right\} \rightarrow_P 1$. To complete the proof, we must show that $\lambda_k(\mathbf{W}^m) = \lambda_k + O_P(m^{-\frac{1}{2}})$.

From the decomposition (A.1) we can write \mathbf{W}^m as a continuous function of:

$$\mathbf{y}_i = \frac{1}{m} \left(k_{i1} \mathbf{u}_i^T, \dots, k_{in} \mathbf{u}_i^T, u_{i1} \mathbf{u}_i^T, \dots, u_{in} \mathbf{u}_i^T \right)^T,$$

and $\frac{1}{m} \mathbf{G}^T \mathbf{\Gamma}^{mT} \mathbf{\Gamma}^m \mathbf{G}$, where $\mathbf{K} = \mathbf{G}^T \mathbf{\Gamma}^{mT}$, this is straightforward for components $(i - iv)$, for component v :

$$\hat{\sigma}_{ave}^2 = \frac{1}{n - \kappa} \sum_{k=\kappa}^n \lambda_k(\mathbf{Z}^m)$$

but $\mathbf{Z}^m = \frac{1}{m} \mathbf{X}^{mT} \mathbf{X}^m$ is a continuous function of \mathbf{y}_i and $\frac{1}{m} \mathbf{G}^T \mathbf{\Gamma}^{mT} \mathbf{\Gamma}^m \mathbf{G}$ and the eigenvalues of \mathbf{Z}^m are a continuous function of the matrix (Henriksen and Isbell, 1953), so $\hat{\sigma}_{ave}^2$ is a continuous function of \mathbf{y}_i .

The expectation of \mathbf{y}_i is $E(\mathbf{y}_i) = (0, \dots, 0, \sigma_i^2, 0, \dots, 0, \sigma_i^2, 0, \dots, 0, \sigma_i^2)^T$. Define $\mathbf{y}_i^* = \sqrt{m} \{\mathbf{y}_i - E(\mathbf{y}_i)\}$; the covariance matrix for this random variable is $\text{cov}(\mathbf{y}_i^*) = \frac{1}{m} \mathbf{\Sigma}_i$. From assumption 1, $\frac{1}{m} \sum_{i=1}^m \mathbf{\Sigma}_i \rightarrow \mathbf{\Sigma}$ and $\sum_{i=1}^m \mathbf{y}_i^*$ is asymptotically normally distributed if the Lindeberg condition holds for every $\epsilon > 0$. Let

$$\psi_i = \left\{ \sum_{j=1}^n (u_{ij}^2 - \sigma_i^2)^2 + \sum_{j=1}^n \sum_{k=1}^n k_{ij}^2 u_{ik}^2 + \sum_{k \neq j} u_{ij}^2 u_{ik}^2 \right\} \mathbf{1}(\|\mathbf{y}_i^*\|^2 > \epsilon).$$

The Lindeberg condition requires $E(\psi_i) \rightarrow 0$ for every i . But ψ_i is only non-zero when:

$$\|\mathbf{y}_i^*\|^2 = \frac{1}{m} \left\{ \sum_{j=1}^n (u_{ij}^2 - \sigma_i^2)^2 + \sum_{j=1}^n \sum_{k=1}^n k_{ij}^2 u_{ik}^2 + \sum_{k \neq j} u_{ij}^2 u_{ik}^2 \right\} > \epsilon$$

an event that has probability zero as $m \rightarrow \infty$, so $\psi_i \rightarrow_P 0$. It is also clear that $|\psi_i| \leq m \|\mathbf{y}_i^*\|^2$ and $E\{m \|\mathbf{y}_i^*\|^2\} < \infty$ by assumption 1. So by the dominated convergence theorem $E(\psi_i) \rightarrow_P 0$ for each i and hence for every $\epsilon > 0$,

$$\sum_{i=1}^m E\{\|\mathbf{y}_i^*\| \mathbf{1}(\|\mathbf{y}_i^*\| > \epsilon)\} = \frac{1}{m} \sum_{i=1}^m E\{\psi_i\} \rightarrow_P 0$$

Since the Lindeberg condition is satisfied $\sum_{i=1}^m \mathbf{y}_i^*$ is asymptotically normally distributed. Since $\text{vec}(\mathbf{W}^m) = \mathbf{g}(\sum_{i=1}^m \mathbf{y}_i) + \text{vec}\left(\frac{1}{m} \mathbf{G} \mathbf{\Gamma}^{mT} \mathbf{\Gamma}^m \mathbf{G}\right)$, where the function $\text{vec}(\cdot)$ concatenates

the columns of a matrix and \mathbf{g} is a continuous function,

$$\sqrt{m} \left(\text{vec}(\mathbf{W}^m) - \text{vec} \left(\frac{1}{m} \mathbf{G}^T \mathbf{\Gamma}^{mT} \mathbf{\Gamma}^m \mathbf{G} \right) \right) \rightarrow \text{MVN}(\mathbf{0}, \mathbf{\Sigma}_w)$$

by the multivariate delta method, so $\sqrt{m}(\mathbf{W}^m - \text{vec}(\frac{1}{m} \mathbf{G}^T \mathbf{\Gamma}^{mT} \mathbf{\Gamma}^m \mathbf{G})) = O_P(1)$. Since $\lambda_r - \lambda_{r+1} = c > 0$ and \mathbf{W}^m is symmetric and real, by Theorem 4.2 of Eaton and Tyler (1991),

$$\begin{aligned} \sqrt{m} \left\{ \lambda_1(\mathbf{W}^m), \dots, \lambda_n(\mathbf{W}^m) \right\}^T - (\lambda_1, \dots, \lambda_n)^T &= O_P(1) \\ \Rightarrow \sqrt{m} \lambda_k(\mathbf{W}^m) &= \sqrt{m} \lambda_k + O_P(1) \quad \forall k. \end{aligned}$$

So $\lambda_k(\mathbf{W}^m) = \lambda_k + O_P(m^{-1/2})$, which completes the proof.

Proof of Corollary 1

$$\begin{aligned} \mathbf{R}^m &= \mathbf{X}^m (\mathbf{I} - \mathbf{S}(\mathbf{S}^T \mathbf{S})^{-1} \mathbf{S}^T) \\ &= \{ \mathbf{B}^m \mathbf{S} + \mathbf{\Gamma}^m \mathbf{G} + \mathbf{U}^m \} \{ \mathbf{I} - \mathbf{S}^T (\mathbf{S} \mathbf{S}^T)^{-1} \mathbf{S} \} \\ &= \mathbf{\Gamma}^m \mathbf{G} + \mathbf{U}^m \mathbf{P}_s \end{aligned}$$

Then we can write:

$$\begin{aligned} \mathbf{W}_R^m &= \frac{1}{m} \mathbf{R}^{mT} \mathbf{R}^m - \hat{\sigma}_{ave}^2 \mathbf{I} \\ &= \underbrace{\frac{1}{m} \mathbf{G}^T \mathbf{\Gamma}^{mT} \mathbf{\Gamma}^m \mathbf{G}}_i + \underbrace{\frac{1}{m} \mathbf{G}^T \mathbf{\Gamma}^{mT} \mathbf{U} \mathbf{P}_s}_ii + \underbrace{\frac{1}{m} \mathbf{P}_s^T \mathbf{U}^{mT} \mathbf{\Gamma}^m \mathbf{G}}_iii + \underbrace{\frac{1}{m} \mathbf{P}_s^T \mathbf{U}^{mT} \mathbf{U}^m \mathbf{P}_s}_iv - \underbrace{\mathbf{P}_s^T \hat{\sigma}_{ave}^2 \mathbf{P}_s}_v \end{aligned}$$

Following the proof of Theorem 1, terms (ii) and (iii) converge to zero, term (i) converges to $\mathbf{G}^T \mathbf{\Delta} \mathbf{G}$, and term (iv) converges to $\mathbf{P}_s^T \bar{\sigma}^2 \mathbf{P}_s = \bar{\sigma}^2 \mathbf{P}_s$. Term (v) is equal to $\hat{\sigma}_{ave}^2 \mathbf{P}_s$, but $\hat{\sigma}_{ave}^2$ is equal to $\frac{1}{n-d-k} \sum_{k=(d+\kappa)}^n \lambda_k(\mathbf{Z}^m)$, which converges almost surely to $\bar{\sigma}^2$. Thus, \mathbf{W}_R^m

converges almost surely to $\mathbf{G}^T \Delta \mathbf{G}$ and the result follows according to the proof of Theorem 1.

Proof of Corollary 2

The proof follows the proof of Lemma 1, where the function, \mathbf{g} , incorporates the project term \mathbf{P}_s .

Web Appendix B

[Figure 1 about here.]

[Figure 2 about here.]

[Figure 3 about here.]

[Figure 4 about here.]

[Figure 5 about here.]

[Figure 6 about here.]

[Figure 7 about here.]

[Figure 8 about here.]

Web Appendix C

[Table 1 about here.]

References

Eaton, M. L. and Tyler, D. E. (1991). On wielandt's inequality and its application to the asymptotic distribution of the eigenvalues of a random symmetric matrix. *Ann Stat* **19**, 260–271.

Feller, W. (1968). *An Introduction to Probability Theory and Its Applications*, volume 1. Wiley, 3 edition.

Henriksen, M. and Isbell, J. R. (1953). On the continuity of the real roots of an algebraic equation. *Proc Amer Math Soc* **4**, 431–4.

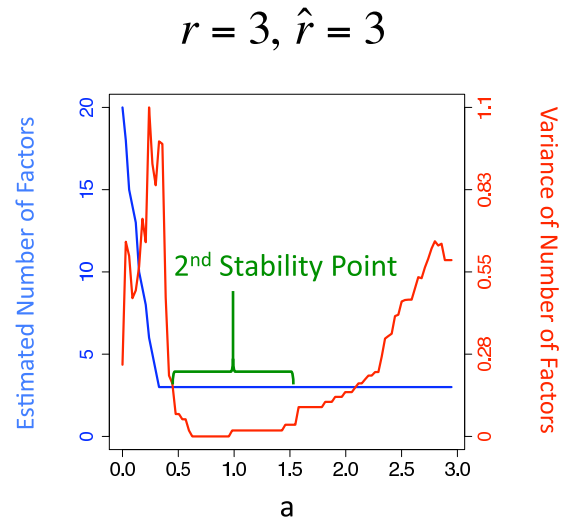


Figure 1. A plot of the estimated number of factors (blue and left axis) and the empirical variance of the estimate for varying set sizes (red and right axis) across a range of coefficients a for a simulated example with $r = 3$. The second stability point (green bracket) is the second point, moving from left to right, where the variance finds a trough. Hallin & Liska (2007) suggest using the estimate corresponding to this second stability point as a practical estimator of the number of factors.

$$r = 3, \hat{r} = 3$$

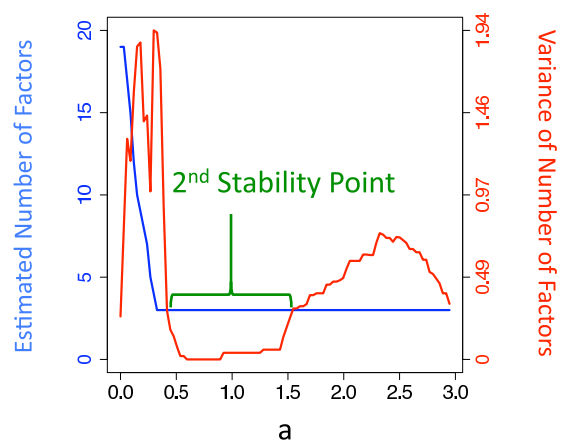


Figure 2. A plot of the estimated number of factors (blue and left axis) and the empirical variance of the estimate for varying set sizes (red and right axis) across a range of coefficients a for a simulated example with $r = 3$. The second stability point (green bracket) is the second point, moving from left to right, where the variance finds a trough. Hallin & Liska (2007) suggest using the estimate corresponding to this second stability point as a practical estimator of the number of factors.

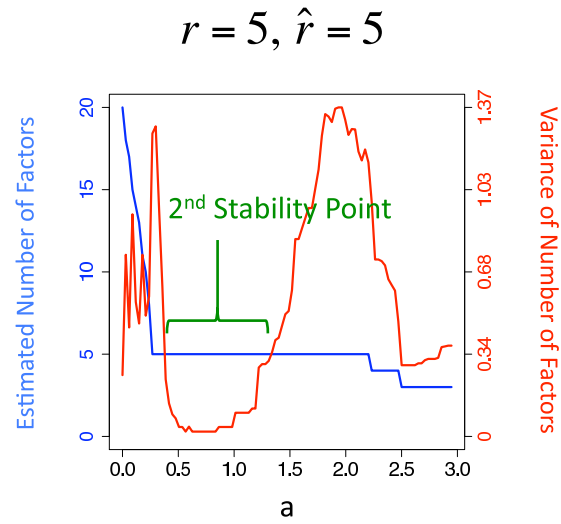


Figure 3. A plot of the estimated number of factors (blue and left axis) and the empirical variance of the estimate for varying set sizes (red and right axis) across a range of coefficients a for a simulated example with $r = 5$. The second stability point (green bracket) is the second point, moving from left to right, where the variance finds a trough. Hallin & Liska (2007) suggest using the estimate corresponding to this second stability point as a practical estimator of the number of factors.

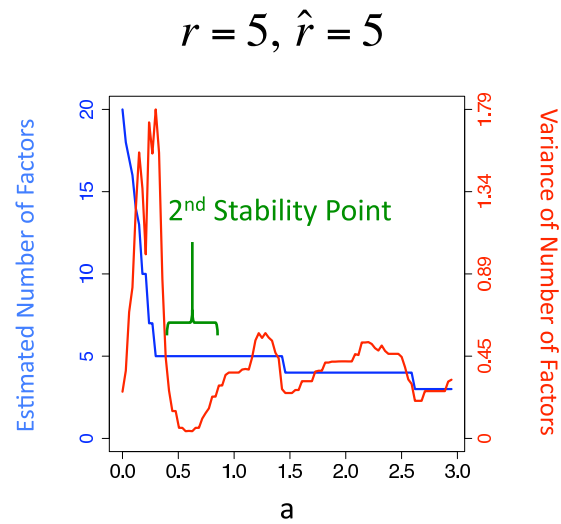


Figure 4. A plot of the estimated number of factors (blue and left axis) and the empirical variance of the estimate for varying set sizes (red and right axis) across a range of coefficients a for a simulated example with $r = 5$. The second stability point (green bracket) is the second point, moving from left to right, where the variance finds a trough. Hallin & Liska (2007) suggest using the estimate corresponding to this second stability point as a practical estimator of the number of factors.

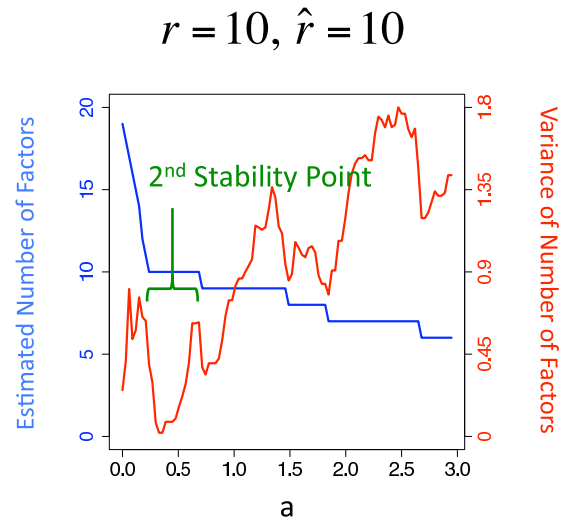


Figure 5. A plot of the estimated number of factors (blue and left axis) and the empirical variance of the estimate for varying set sizes (red and right axis) across a range of coefficients a for a simulated example with $r = 10$. The second stability point (green bracket) is the second point, moving from left to right, where the variance finds a trough. Hallin & Liska (2007) suggest using the estimate corresponding to this second stability point as a practical estimator of the number of factors.

$$r = 10, \hat{r} = 10$$

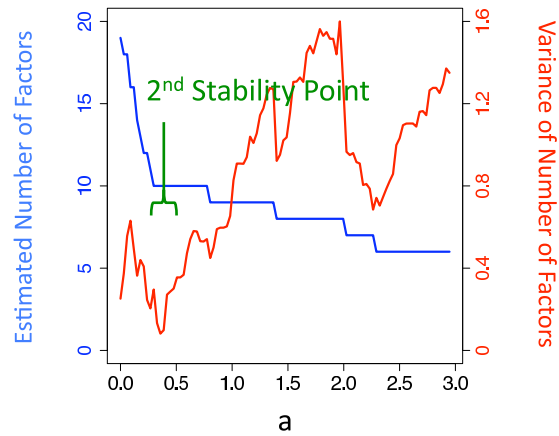


Figure 6. A plot of the estimated number of factors (blue and left axis) and the empirical variance of the estimate for varying set sizes (red and right axis) across a range of coefficients a for a simulated example with $r = 10$. The second stability point (green bracket) is the second point, moving from left to right, where the variance finds a trough. Hallin & Liska (2007) suggest using the estimate corresponding to this second stability point as a practical estimator of the number of factors.

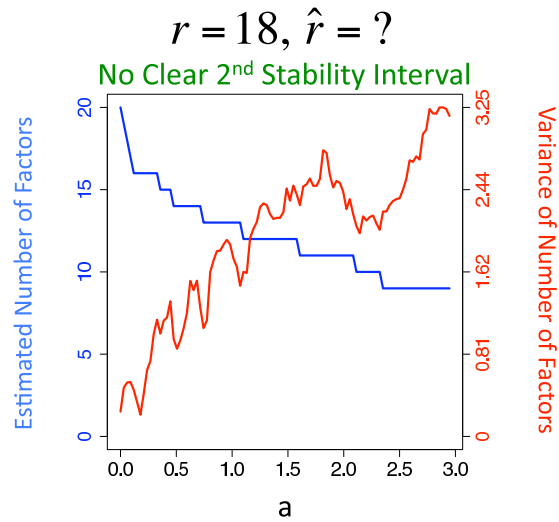


Figure 7. A plot of the estimated number of factors (blue and left axis) and the empirical variance of the estimate for varying set sizes (red and right axis) across a range of coefficients a for a simulated example with $r = 18$. The second stability point (green bracket) is the second point, moving from left to right, where the variance finds a trough. Hallin & Liska (2007) suggest using the estimate corresponding to this second stability point as a practical estimator of the number of factors. Since $r = 18$ is close to the sample size $n = 20$, there is no clear second stability point, so the Hallin & Liska approach does not give an estimate.

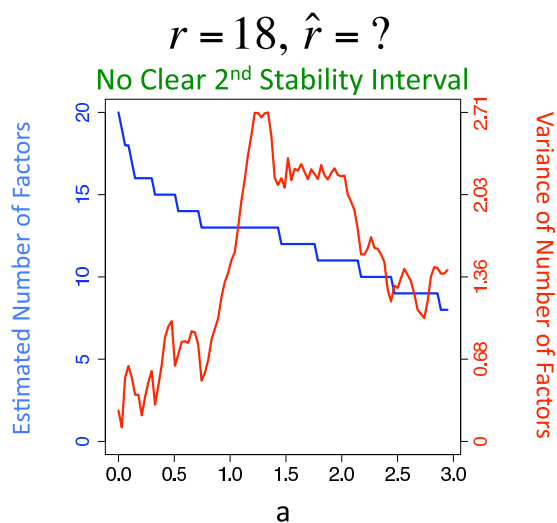


Figure 8. A plot of the estimated number of factors (blue and left axis) and the empirical variance of the estimate for varying set sizes (red and right axis) across a range of coefficients a for a simulated example with $r = 18$. The second stability point (green bracket) is the second point, moving from left to right, where the variance finds a trough. Hallin & Liska (2007) suggest using the estimate corresponding to this second stability point as a practical estimator of the number of factors. Since $r = 18$ is close to the sample size $n = 20$, there is no clear second stability point, so the Hallin & Liska approach does not give an estimate.

Table 1

Results from a simulation experiment. For each combination of m, n and r , 100 independent microarray data sets were simulated. The average (s.d.) RMSFE, a measure of how well the eigenvectors of W_m span the linear space spanned by G , is reported for the Lemma 1 estimator of r and the Bai & Ng (2002) and Buja & Eyuboglu (1992) estimators.

(m, n)	r	RMSFE(\hat{r}) $\times 10^5$	RMSFE(\hat{r}_{bn}) $\times 10^5$	RMSFE(\hat{r}_{be}) $\times 10^5$
(1000,10)	3	403.84 (773.88)	2514.19 (2881.21)	915.18 (1571.98)
(5000,10)	3	78.44 (271.25)	3022.17 (2409.22)	785.95 (1446.20)
(10000,10)	3	50.18 (215.60)	2881.82 (2585.60)	817.71 (1753.53)
(1000,20)	3	109.77 (23.18)	685.79 (1683.42)	133.18 (237.64)
(5000,20)	3	22.49 (4.67)	426.71 (1266.14)	22.49 (4.67)
(10000,20)	3	10.42 (2.22)	193.79 (972.41)	10.42 (2.22)
(1000,100)	3	101.24 (7.29)	101.24 (7.29)	101.24 (7.29)
(5000,100)	3	20.26 (1.57)	20.26 (1.57)	20.26 (1.57)
(10000,100)	3	10.11(0.82)	10.11(0.82)	10.11 (0.82)
(1000,10)	5	822.46 (624.17)	2514 .19 (1704.76)	3256.58 (1698.78)
(5000,10)	5	224.48(309.48)	3022.17 (1729.20)	2885.43 (1636.53)
(10000,10)	5	129.70(245.14)	2606.50 (1371.87)	2655.60 (1474.46)
(1000,20)	5	395.34 (591.32)	1378.21 (1101.71)	480.67 (718.98)
(5000,20)	5	46.95 (151.17)	1298.78 (1109.70)	230.22 (563.62)
(10000,20)	5	18.00 (78.65)	980.16 (991.06)	191.60 (478.68)
(1000,100)	5	99.65 (7.22)	99.65 (7.22)	99.65 (7.21)
(5000,100)	5	20.11(1.31)	20.11 (1.31)	20.11 (1.31)
(10000,100)	5	10.10 (0.63)	10.10 (0.63)	10.10 (0.63)
(1000,20)	10	1108.28 (435.75)	1732.07 (644.47)	3034.82 (1076.86)
(5000,20)	10	297.95 (245.09)	1715.14 (545.61)	2683.19 (717.63)
(10000,20)	10	198.12 (164.99)	1592.15 (520.36)	2562.26 (711.42)
(1000,100)	10	96.40 (5.31)	96.40 (5.31)	96.40 (5.31)
(5000,100)	10	19.30 (1.06)	19.30 (1.06)	19.30 (1.06)
(10000,100)	10	9.60 (0.44)	9.60(0.44)	9.60 (0.44)
(1000,20)	18	1148.48 (273.08)	1562.41(329.37)	5479.46 (1011.77)
(5000,20)	18	497.17(144.31)	1484.74 (312.31)	5033.88 (832.27)
(10000,20)	18	336.51 (104.66)	1472.69 (350.24)	4876.27 (867.75)
(1000,100)	18	314.30 (293.42)	90.25 (3.61)	97.64 (53.06)
(5000,100)	18	17.88 (0.74)	17.88 (0.74)	17.88 (0.74)
(10000,100)	18	8.88 (0.34)	8.88 (0.34)	8.88 (0.34)