

Gently Clarifying the Application of Horn's Parallel Analysis to Principal Component Analysis Versus Factor Analysis

Alexis Dinno
Portland State University

March 18, 2014

Introduction

Horn's parallel analysis (PA) is an empirical method for deciding how many components in a principal component analysis (PCA) or factors in a common factor analysis (CFA) drive the variance observed in a data set of n observations on p variables (Horn, 1965). This decision of how many components or factors to retain is critical in applications of PCA or CFA to reducing the dimensionality of data in analysis (as when compositing multiple scale items into a single score), and also in exploratory factor analysis where the different contributions of each factor to each observed variable help generate theory (Preacher & MacCallum, 2003; Velicer & Jackson, 1990). As will be shown, the development of PA was predicated upon properties of PCA. However, some have been exponents of the use of PA for CFA (Velicer, Eaton, & Fava, 2000). The correct application of PA with CFA requires modification to the original PA procedure. This paper attempts to clarify PA with respect to both PCA and CFA.

Concerning eigenvalues in PCA and CFA

PCA and CFA are two similar methods describing complex correlation in an n by p matrix \mathbf{X} of observed data. Both methods produce eigenvalues— λ s ordered in magnitude from largest (λ_1) to smallest (λ_p)—which apportion variance along p unobserved dimensions. One major interpretive difference between PCA and CFA, is that in the former, each (unrotated) eigenvalue represents a portion of total variance in \mathbf{X} , and in the later each (unrotated) eigenvalue represents a portion of common standardized variance shared among all p variables. This means that in PCA $\Sigma(\mathbf{\Lambda}) = p$, and that in CFA $\Sigma(\mathbf{\Lambda}) < p$ (also: eigenvalues from a CFA can be negative).

Assume that PCA is taken to be a function of observed n by p data \mathbf{X} that returns a set of p eigenvalues. If $e(\mathbf{A})$ is a function returning a vector of eigenvalues of square matrix \mathbf{A} , and $\text{cor}(\mathbf{X})$ is the correlation matrix of \mathbf{X} , then, leaving out the issue of eigenvectors, a PCA of \mathbf{X} returns the vector $\mathbf{\Lambda}$ of eigenvalues as in (1). For PCA of $\text{cor}(\mathbf{X})$, each observed variable is assumed to contribute an equal portion to total variance.

$$\mathbf{\Lambda}_{\mathbf{X}} = e(\text{cor}(\mathbf{X})) \quad (1)$$

Where

$$\mathbf{\Lambda}_{\mathbf{X}} = [\lambda_1, \lambda_2, \dots, \lambda_p] \quad (2)$$

and $\lambda_1 > \lambda_2 > \dots > \lambda_p$.

If \mathbf{U} is a matrix of n observations of p uncorrelated variables, then as n approaches ∞ , $\mathbf{\Lambda}_{\mathbf{U}}$ approaches the 1 by p unit vector $\mathbf{1}$ (3). This jibes with the substantive interpretation of PCA as apportioning total standardized variance: if p variables are perfectly uncorrelated, then in an infinite population they must each explain exactly the same amount of standardized variance, namely $(1/p) \times p$, or 1.

$$\lim_{n \rightarrow \infty} \mathbf{\Lambda}_{\mathbf{U}} = \mathbf{1}_{1 \times p} \quad (3)$$

One can easily demonstrate this limiting property by running the series of commands in R listed in Appendix A which return the eigenvalues of a PCA of $\text{cor}(\mathbf{U})$ for progressively larger values of n for $p = 20$.

If $\text{cov}(\mathbf{X})$ is the variance-covariance matrix of \mathbf{X} , then, again leaving out the issue of eigenvectors, a PCA of \mathbf{X} could, instead return the vector $\mathbf{\Lambda}$ of eigenvalues as in (4). Here the interpretation of the eigenvalues $\mathbf{\Lambda}$ is altered, so that the i^{th} eigenvalue explains $\lambda_i / \text{trace}(\text{cov}(\mathbf{X}))$ proportion of the total variance, $\text{trace}(\text{cov}(\mathbf{X}))$ is the sum of the diagonal of the variance-covariance matrix, and each observed variable is assumed to contribute the absolute magnitude of its variance to total variance. Verify (3) for a PCA of $\text{cov}(\mathbf{U})$ using the R commands in Appendix B.

$$\mathbf{\Lambda}_{\mathbf{X}} = e(\text{cov}(\mathbf{X})) \quad (4)$$

Alexis Dinno, School of Community Health, Portland State University.

Correspondence concerning this paper should be addressed to Alexis Dinno, School of Community Health, Portland State University, PO Box 751-SCH, Portland, Oregon 97207-0751 United States of America. Email alexis.dinno@pdx.edu

This paper should be cited as an unpublished manuscript using the URL: http://doyenne.com/Software/files/PA_for_PCA_vs_FA.pdf

The behavior of CFA relevant to PA in the limit of n can be approached in the same fashion. If the function $\text{diag}(\mathbf{A})$ of a square matrix returns a square matrix with the main diagonal elements (a_{ij} where $i = j$) of \mathbf{A} , and zeros in all other elements, and if \mathbf{A}^+ is the Moore-Penrose inverse (also ‘generalized inverse’, or ‘pseudoinverse’) of the matrix \mathbf{A} , then a CFA of \mathbf{X} returns the vector $\Lambda_{\mathbf{X}}$ of eigenvalues as in (5).

$$\Lambda_{\mathbf{X}} = e\left(\text{cor}(\mathbf{X}) - \text{diag}(\text{cor}(\mathbf{X})^+)\right) \quad (5)$$

and $\Lambda_{\mathbf{X}} = [\lambda_1, \lambda_2, \dots, \lambda_p]$, with $\lambda_1 > \lambda_2 > \dots > \lambda_p$ as in (2).

If \mathbf{U} is a matrix of n observations on p uncorrelated variables, then as n approaches ∞ , $\Lambda_{\mathbf{U}}$ approaches the 1 by p zero vector $\mathbf{0}$ (6). This jibes with the substantive interpretation of common factor analysis as apportioning common standardized variance: if p variables are perfectly uncorrelated, then in an infinite population there can be no common standardized variance, so each factor ‘explains’ zero common variance.

$$\lim_{n \rightarrow \infty} \Lambda_{\mathbf{U}} = \mathbf{0}_{1 \times p} \quad (6)$$

One can easily demonstrate this property by running the series of commands in R listed in Appendix C (requires the MASS package from <http://cran.r-project.org>) which return the eigenvalues of \mathbf{U} for progressively larger values of n (the commands return the diagonal of $\Lambda_{\mathbf{U}}$).

The difference between (3) and (6) is critical to the correct application of PA to PCA versus CFA.

Applying PA

Kaiser (1960) asserted that in application of PCA one would retain those components with eigenvalues greater than one (7).

$$\lambda_q \begin{cases} > 1 & \text{retain} \\ \leq 1 & \text{do not retain} \end{cases} \quad (7)$$

Where q indexes the eigenvalues from 1 to p .

Horn (1965) elaborated upon this logic by pointing out that applied researchers do not have an infinite number of observations. According to Horn, in order to account for ‘sampling error and least squares bias’ due to finite n , one would want to:

1. conduct a parallel PCA on an n by p matrix of uncorrelated random values;
2. repeat this k times;
3. average each of the eigenvalues λ_q^r over k , to produce $\bar{\lambda}_q^r$; and
4. adjust λ_q by subtracting from it $(\bar{\lambda}_q^r - 1)$ to produce λ_q^{adj} .

The retention criterion of PA is to retain those first components with adjusted eigenvalues greater than one (8). Technically, PA is a stopping rule in PCA, because the adjustment to subsequent components—especially the last few components—may sometimes increase their eigenvalues above the value of one. The retention criterion in (8) can be stated in a mathematically equivalent way as ‘retain those first components with unadjusted eigenvalues greater than the corresponding mean eigenvalue of random data’ (9).

$$\lambda_q^{adj} \begin{cases} > 1 & \text{retain} \\ \leq 1 & \text{do not retain (and stop)} \end{cases} \quad (8)$$

$$\lambda_q \begin{cases} > \bar{\lambda}_q^r & \text{retain} \\ \leq \bar{\lambda}_q^r & \text{do not retain (and stop)} \end{cases} \quad (9)$$

PA must be amended for use with CFA by calculating the adjusted eigenvalue λ_q^{adj} as $\lambda_q - \bar{\lambda}_q^r$. The retention criteria must likewise be changed to retain those first adjusted eigenvalues greater than zero (10). Technically, PA is a stopping rule in CFA, because the adjustment to subsequent common factors—especially the last few factors—may sometimes increase their eigenvalues above the value of one. And as with PA for PCA, PA for CFA can be restated in an equivalent form as ‘retain those unadjusted eigenvalues greater than the corresponding mean eigenvalue of random data’ (9).

$$\lambda_q^{adj} \begin{cases} > 0 & \text{retain} \\ \leq 0 & \text{do not retain (and stop)} \end{cases} \quad (10)$$

References

- Gorsuch, R. L. (1983). *Factor Analysis 2nd ed.* NJ: Lawrence Erlbaum Associates.
- Horn, J. (1965). A rationale and test for the number of factors in factor analysis. *Psychometrika*, 30(2), 179–185.
- Kaiser, H. (1960). The Application of Electronic Computers to Factor Analysis. *Educational and Psychological Measurement*, 20(1), 141–151.
- Preacher, K. J. & MacCallum, R. C. (2003). Repairing Tom Swift’s Electric Factor Analysis Machine. *Understanding Statistics*, 2(1), 13–43.
- Velicer, W. F., Eaton, C. A., & Fava, J. L. (2000). Construct explication through factor or component analysis: A review and evaluation of alternative procedures for determining the number of factors or components. In R. D. Goffen & E. Helms (Eds.), *Problems and Solutions in Human Assessment - Honoring Douglas N. Jackson at Seventy* (pp. 41–71). Springer.
- Velicer, W. F. & Jackson, D. N. (1990). Component analysis versus common factor analysis: Some further observations. *Multivariate Behavioral Research*, 25(1), 97–114.

NOTE: Both the `verbatim` and `listings` approaches to representing code produce output that has problems for simple select-cut-paste operations with \LaTeX on my computer. This means that selecting, copying, pasting and then trying to execute the following examples may present difficulties, so I recommend typing them out. Email me if you have suggestions for how to fix this.

Appendix A

*

The limiting case of Λ in PCAs of the correlation matrix of uncorrelated data

```
p <- 20
for (n in c(100, 1000, 1000000)) {
  U <- matrix(rnorm(n*p),n,p)
  Lambda_U <- eigen(cor(U), only.values = TRUE)[[1]]
  cat("For n = ", n, ", Lambda_U (PCA) = \n", sep="")
  print(Lambda_U)
  cat("\n")
}
```

Appendix B

*

The limiting case of Λ in PCAs of the variance-covariance matrix of uncorrelated data with differing variances

```
p <- 20
for (n in c(100, 1000, 1000000)) {
  U <- matrix(rnorm(n*p, sd=sqrt(runif(n=p, min=1, max=50))),n,p, byrow=TRUE)
  Lambda_U <- eigen(cor(U), only.values = TRUE)[[1]]
  cat("For n = ", n, ", Lambda_U (PCA) = \n", sep="")
  print(Lambda_U)
  cat("\n")
}
```

Appendix C

*

The limiting case of Λ in CFAs of uncorrelated data

```
library(MASS)
p <- 20
for (n in c(100, 1000, 1000000)) {
  U <- matrix(rnorm(n*p),n,p)
  eigen(cor(U)-ginv(diag(diag(ginv(cor(U))))), only.values = TRUE)[[1]]
  cat("For n = ", n, ", Lambda_U (CFA) = \n", sep="")
  print(Lambda_U)
  cat("\n")
}
```