# Standardization and Estimation of the Number of Factors for Panel Data[*]

Ryan Greenaway-McGrevy[†]    Chirok Han[‡]    Donggyu Sul[§]

**Abstract**    Practitioners often standardize panel data before estimating a factor model. In this paper we show an example that the standardization leads to inconsistent estimation of the factor number. When the common component exhibits strong heteroskedasticity, the conventional eigenvalue-based decompositions are consistent but standardization does not necessarily result in consistent estimation. To overcome this issue, we recommend using a "minimum-rule" whereby the minimum factor-number estimated from both the conventional and standardized panel is used. Monte Carlo studies and an empirical application are provided.

**Keywords**    Factor Model, Selection Criteria, Principal Components Estimator, Bai-Ng Criteria, Standarization, Panel Data

**JEL Classification**    C33

---

## 1. INTRODUCTION

Factor models are increasingly being used in empirical economics. They are used in forecasting (Stock and Watson, 2002), macroeconomic modelling and policy analysis (Bernanke and Boivin, 2003; Bernanke et al., 2005), and more recently in price index construction (Bajari and Benkard, 2005; Oulton, 2008).

Approximate factor models permit weak heteroskedasticity (e.g. Chamberlain and Rothschild, 1983; Bai and Ng, 2002a). Under weak heteroskedasticity the principal component (PC) decomposition can be used to consistently estimate the factor-number and factor structure as $N$ (cross sections) and $T$ (time series) grow large (e.g. Bai, 2003; Bai and Ng, 2002a). In practice however the heteroskedasticity may be so prominent that the assumption of weak heteroskedasticity is inappropriate. Conventional solution to this problem is standardization: Standardize the data by dividing each time-series in the panel by its sample standard deviation. For example, Stock and Watson (2002), Bai and Ng (2006, 2008), Boivin and Ng (2006), and Kapetanios and Pesaran (2007) standardize the data in their respective empirical sections prior to running a PC decomposition. (In addition, statistical packages such as Stata and EViews offer principal component decompositions based on standardized data.) Standardization ensures that each re-scaled series has unit sample variance, such that by construction no heteroskedasticity (in terms of within-group sample variance) remains in the observed series. Greenaway-McGrevy *et al.* (2012b) show that the factor number estimators based on the PC decomposition of standardized data are consistent when the source of heteroskedasticity is the idiosynratic component.

The main purpose of the present paper is to examine a potential problem of standardization with regard to the estimation of the number of the common factors. We find that the validity of standardization is dependent on whether the main source of heteroskedasticity is the idiosyncratic component or the common component. When heteroskedasticity mostly comes from the idiosyncratic error, the PC estimation of the factor model would be consistent when applied to the standardized panel. But if the panel has excessive heteroskedasticity due to unbounded factor loadings, then standardization may create heteroskedasticity in the re-scaled idiosyncratic error and lead to inconsistency. In particular, the IC criteria (Bai and Ng, 2002a) applied to the standardized panel may over-estimate the factor number.

Yet in practice the source of the heteroskedasticity is not known, so the econometrician may not be able to tell whether standardization will help or hinder identification of the factor model. Our findings do however provide some

advice for practitioners to overcome this problem. Noting that estimators over-estimate the factor-number when inconsistent, we use a minimum rule (as done by Greenaway-McGrevy et al., 2012a, for panels with strong serial correlation), whereby the minimum factor-number estimated from both the conventional and standardized panel is used as a final estimate. This minimum rule would restore consistency of the factor number estimation in many cases when the source of heteroskedasticity is unclear.

The remainder of the paper is organized as follows. In the next section we discuss the cases where the standardization method fails. Section 3 reports results of Monte-Carlo studies. Section 4 presents an empirical example, in which the Bai-Ng $IC_{p2}$ criterion gives a reasonable factor number without standardization but standardization yields unreasonable estimates. Section 5 concludes.

## 2. FAILURE OF STANDARDIZATION

Consider the factor model $X = F\Lambda' + e$ in matrix notation, where $X$ and $e$ are the $T \times N$ matrices of $X_{it}$ and $e_{it}$ respectively, $F$ is the $T \times r$ matrix of common factors, and $\Lambda$ is the $N \times r$ matrix of factor loadings. We are interested in determining the rank $r$ of important common shocks $F$ using the data $X$.

Let $k_{\max}$ be given. The IC criteria of Bai and Ng (2002a) are defined by $IC(k) = \ln V(k) + kg(N,T)$ for some $g(N,T)$, where $V(k)$ is the minimum value of $\sum_{i=1}^{N}\sum_{t=1}^{T}(X_{it} - \lambda_i^{k'}F_t^k)^2$ that can be attained by adjusting $\lambda_i^k \in \mathbb{R}^k$ and $F_t^k \in \mathbb{R}^k$ for $i = 1, \ldots, N$ and $t = 1, \ldots, T$. Then the minimizer $\hat{k}$ of $IC(k)$ over $k = 0, 1, \ldots, k_{\max}$ is consistent under weak heteroskedasticity and serial correlation for idiosyncratic errors and other regularity conditions (Bai and Ng, 2002a). Greenaway *et al.* (2012b) provide a more flexible proof of consistency which can be extended to the case of standardization. The part of their results relevant to the current subject is rephrased below as Theorem 1. Let $\text{eigval}_j(\cdot)$ stand for the $j$th largest eigenvalue of the argument. The minimal eigenvalue is also denoted by $\text{eigval}_{\min}(\cdot)$. We denote by $\|A\|$ the square-root of the largest eigenvalue of $A'A$, i.e., the largest singular value of $A$, which is also called the matrix 2-norm of $A$.

**Theorem 1 (from Greenaway-McGrevy et al., 2012b, Theorem 1).** *Suppose that $T^{-1}F'F$ converges to a finite nonsingular matrix. If (i) $\text{eigval}_{\min}(N^{-1}\Lambda'\Lambda) \geq m > 0$ for all $N$, (ii) $\|\frac{1}{NT}ee'\| = O_p(C_{NT}^{-2})$, where $C_{NT} = \min(N^{1/2}, T^{1/2})$, and (iii) $\text{eigval}_1(\frac{1}{NT}ee')/\text{eigval}_L(\frac{1}{NT}ee') = O_p(1)$ for some $L$ such that $Lg(N,T) \to \infty$, then $P(\hat{k} = r) \to 1$ as $N, T \to \infty$.*

In the above result, the maintained supposition that $T^{-1}F'F$ converges to a finite nonsingular matrix is standard in the literature, but it is notable that integrated factors are not allowed. Extension to I(1) factors would be straightforward if wanted. Condition (i) is introduced to exclude the case of zero or very weak factor loadings. Condition (ii) loosely means that heteroskedasticity in the idiosyncratic error is limited within a sufficiently small boundary for the given factor loadings. For example, if $e_{it}$ is *iid* over $i$ and $t$, then this condition is satisfied and conditions for more general cases with weak heteroskedasticity and serial correlation can be derived using the results of Yin *et al.* (1988). (See Bai and Ng, 2002b.) Condition (iii) is important but has been overlooked in the literature, and states that a sufficiently large number of cross sectional units have sufficiently homoskedastic idiosyncratic errors *relative to* other idiosyncratic error series. As a counter example, if $Ee_{it}^2 = 1$ for $i = 1, 2$ and $Ee_{it}^2 = 0$ for all other $i$, then Condition (ii) is satisfied but Condition (iii) is violated. More detailed remarks are found in Greenaway-McGrevy *et al.* (2012b).

It is notable that Theorem 1 allows for large heteroskedasticity in the common component due to large $\lambda_i$ and, loosely speaking, requires that there are no relatively prominent idiosyncratic errors which can be mistaken for common factors.

Theorem 1 is more flexible than the usual proofs in the literature as the conditions contain no population expectations, and is useful for investigating the asymptotic behavior of factor number estimates after standardization (c.f., Bai and Ng, 2002a). In particular, Theorem 1 allows us to check the consistency of factor number estimates based on standardized data $X_{it}/\hat{\sigma}_{Xi}$, where $\hat{\sigma}_{Xi}^2 = \frac{1}{T-1}\sum_{t=1}^{T}(X_{it} - \bar{X}_i)^2$ and $\bar{X}_i = \frac{1}{T}\sum_{t=1}^{T} X_{it}$. Specifically, as $\hat{\sigma}_{Xi}^{-1}X_{it} = (\hat{\sigma}_{Xi}^{-1}\lambda_i)'F_t + \hat{\sigma}_{Xi}^{-1}e_{it}$, we can see that if the rescaled factor loadings $\hat{\sigma}_{Xi}^{-1}\lambda_i$ and the rescaled idiosyncratic errors $\hat{\sigma}_{Xi}^{-1}e_{it}$ satisfy the conditions for Theorem 1, then the factor number estimator $\hat{k}_{std}$ after standardization is consistent.

With regard to heteroskedasticity, we pay attention to Condition (iii) of Theorem 1. Especially, if $\hat{\sigma}_{Xi}^2$ is large due to large $\lambda_i$ for some $i$ and $Ee_{it}^2$ is not proportionally large, then the volatility of $\hat{\sigma}_{Xi}^{-1}e_{it}$ will be small relative to other cross sectional units. This can lead to the violation of Condition (iii) of Theorem 1, and $\hat{k}_{std}$ may be larger than the true factor number with nontrivial probabilities. The following example illustrates this possibility.

**Example 1.** Let $X_{it} = c_i\tilde{\lambda}_iF_t + e_{it}$, where the scalar $\tilde{\lambda}_i$ are bounded, $\tilde{\lambda}_i^2$ are uniformly sufficiently bounded away from zero, $N^{-1}\sum_{i=1}^{N}\tilde{\lambda}_i^2 \to 1$, $e_{it} \sim iid\ N(0,1)$, $F_t \sim iid\ N(0,I_r)$, and $e_{it}$, $F_t$, and $\tilde{\lambda}_i$ are mutually independent. In short, $\tilde{\lambda}_i$, $F_t$ and $e_{it}$ behave regularly. Now let $c_i = 1$ for $i = 1$, and $c_i = \sqrt{N}$ for $i \geq 1$. The true

factor number is 1. All the conditions for consistent estimation (Theorem 1) are satisfied in this setting, and the usual factor number estimator is consistent.

Standardization, on the other hand, leads to a different result. We have $\hat{\sigma}_{Xi}^2 = (c_i^2 \tilde{\lambda}_i^2 + 1)[1 + o_p(1)]$, and the standardized data are

$$\frac{X_{it}}{\hat{\sigma}_{Xi}} = \left( \frac{c_i \tilde{\lambda}_i}{\hat{\sigma}_{Xi}} \right) F_t + \frac{e_{it}}{\hat{\sigma}_{Xi}}.$$

Because $\hat{\sigma}_{Xi}/c_i$ is stochastically bounded asymptotically, the rescaled factor loadings $(c_i/\hat{\sigma}_{Xi})\tilde{\lambda}_i$ do not shrink to zero, and Condition (i) of Theorem 1 would be satisfied by the standardized factor loadings. Condition (ii) of Theorem 1 is also satisfied because $1/\hat{\sigma}_{Xi}$ is bounded (asymptotically). However, the standardized idiosyncratic errors $\hat{\sigma}_{Xi}^{-1} e_{it}$ show more temporal variability for $i = 1$ than for all other $i$, because $\hat{\sigma}_{X1}^{-1} = O_p(1)$ and $\hat{\sigma}_{Xi}^{-1} = O_p(N^{-1/2})$ for $i > 1$. That is, the first cross-sectional unit shows a *relatively* larger idiosyncratic variation compared to the rest, and Condition (iii) of Theorem 1 is violated by $\hat{\sigma}_{Xi}^{-1} e_{it}$ and the factor number estimator after standardization is inconsistent.[1] This is illustrated by a Monte Carlo study in the following section. □

Example 1 shows that heterogeneity in factor loadings can be converted to heteroskedasticity in idiosyncratic errors by standardization. In this example, heteroskedasticity is present due to the factor loadings, so the usual PC estimation (without standardization) is consistent as shown by Theorem 1, but standardization may lead to inconsistency. Table 1 summarizes the relationship between the source of heteroskedasticity and the result of standardization. Note that standardization always works if there are no common factors in the panel data because then the source of heteroskedasticity is always the idiosyncratic errors.

Table 1: Source of heteroskedasticity and effect of standardization

|  | Source of heteroskedasticity | |
|---|---|---|
|  | $\lambda_i$ | $e_{it}$ |
| No standardization | Consistent | Possible over-estimation |
| Standardization | Possible over-estimation | Consistent |

[1] As the conditions in Theorem 1 are sufficient but not necessary, the violation of Condition (iii) in Example 1 does not necessarily imply inconsistency of $\hat{k}_{std}$. But we can show that the factor number is over-estimated for this example.

Because the components $\lambda_i' F_t$ and $e_{it}$ of $X_{it}$ are not separately observed, it is hard to identify the source of excess heteroskedasticity. Also standardization does not always give a satisfactory solution to the over-estimation problem as Example 1 and Table 1 show. To make things worse, if the sources of heteroskedasticity are different across $i$, then it would be even harder (if even possible) to tell whether an estimate is consistent or not.

However, regardless of whether the excessive volatility is resultant from either the factor loading or the idiosyncratic error, heteroskedasticity usually (though not always) results in overestimation rather than underestimation. Thus, as in Greenaway-McGrevy *et al.* (2012a), we may apply a minimum rule given by

$$\hat{k}_{\min} = \min\left[\hat{k}_{\text{std}}, \hat{k}_{\text{no-std}}\right] \tag{1}$$

to estimate the factor number. Here $\hat{k}_{\text{std}}$ and $\hat{k}_{\text{no-std}}$ are the factor number estimate with and without standardization, respectively. If $\hat{k}_{\text{std}} < \hat{k}_{\text{no-std}}$, it would imply that the source of heteroskedasticity is idiosyncratic errors, and vice versa. Common factors may be estimated using the standardized data or the raw data in accordance if the raw data or the standardized data gives a consistent factor number estimator.

## 3. MONTE CARLO STUDIES

In this section, we verify our theoretical claim by means of Monte Carlo experiments based on two simple data generating processes. We generate $X_{it}$ according to $X_{it} = \lambda_i' F_t + e_{it}$ for the single factor ($r = 1$). We consider two cases:

Case 1:   $\sigma_{ei}^2 = \begin{cases} N^{1/2} \text{ for } i = 1 \\ \phantom{N^{1/2}} 1 \text{ for } i \geq 2 \end{cases}$   $\lambda_i \sim iid\ N(0,1)$ for all $i$

Case 2:   $\sigma_{ei}^2 = 1$ for all $i$   $\lambda_i = \begin{cases} 0 \text{ for } i = 1 \\ N^{1/2} - 1 \text{ for } i \geq 2. \end{cases}$

and where $e_{it} \sim iid\ N\left(0, \sigma_{ei}^2\right)$ and $F_t \sim iid\ N(0,1)$. Case 1 investigates heteroskedasticity in the idiosyncratic component. The true number of common factor $r$ is equal to one, and eigenvalue-based criteria will asymptotically select two factors when applied directly to $X_{it}$. Case 2 investigates heteroskedasticity in the factor loading coefficients. We consider $T, N = 25, 50, 100,$. We use the $IC_{p2}$ criterion of Bai and Ng (2002a), which is $IC_{p2}(k) = \ln V(k) + k(N^{-1} + T^{-1}) \ln C_{NT}^2$. Simulations are replicated 10,000 times, and we set $k_{\max}$ to 3.

Table 2 shows the selection percentage before and after standardization for Case 1. (Reported figures are rounded to the nearest percent.) As expected,

Table 2: Heteroskedastic idiosyncratic errors (Case 1) with $r = 1$ (percentage of selection).

| | | Without standardization | | | | With standardization | | | | Minimum rule | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | $N$ | k=0 | k=1 | k=2 | k=3 | k=0 | k=1 | k=2 | k=3 | k=0 | k=1 | k=2 | k=3 |
| 25 | 25 | 0 | 63 | 37 | 0 | 1 | 99 | 0 | 0 | 1 | 99 | 0 | 0 |
| 50 | 25 | 0 | 27 | 72 | 0 | 0 | 100 | 0 | 0 | 0 | 100 | 0 | 0 |
| 100 | 25 | 0 | 5 | 95 | 0 | 0 | 100 | 0 | 0 | 0 | 100 | 0 | 0 |
| 25 | 50 | 0 | 61 | 39 | 0 | 0 | 100 | 0 | 0 | 0 | 100 | 0 | 0 |
| 50 | 50 | 0 | 48 | 52 | 0 | 0 | 100 | 0 | 0 | 0 | 100 | 0 | 0 |
| 100 | 50 | 0 | 5 | 95 | 0 | 0 | 100 | 0 | 0 | 0 | 100 | 0 | 0 |
| 25 | 100 | 0 | 75 | 25 | 0 | 0 | 100 | 0 | 0 | 0 | 100 | 0 | 0 |
| 50 | 100 | 0 | 49 | 51 | 0 | 0 | 100 | 0 | 0 | 0 | 100 | 0 | 0 |
| 100 | 100 | 0 | 12 | 88 | 0 | 0 | 100 | 0 | 0 | 0 | 100 | 0 | 0 |

the $\mathrm{IC}_{p2}(k)$ selection criterion without standardization ove-restimates the factor number. As $N$ and $T$ increase $\mathrm{IC}_{p2}(k)$ consistently selects two common factors rather than the true $r = 1$. However after standardization the $\mathrm{IC}_{p2}(k)$ criteria is consistent.

Table 3 reports the results for Case 2, where heteroskedasticity originates from factor loadings and standardization fails asymptotically. (Reported figures are rounded to the nearest percent.) However, without standardization the factor number is consistently estimated.

In both Case 1 and Case 2, the minimum rule provides a consistent estimator for the factor number.

## 4. EMPIRICAL EXAMPLE

We have discussed that standardization can lead to over-estimation of the true factor number when strong heteroskedasticity is driven by the common component. We now consider an empirical example in which standardization indeed leads to a larger factor number estimate than that obtained without standardization. We estimate the number of common factors to annual growth rates in US value-added by industry. We use the Bureau of Economic Analysis' (BEA) annual quantity indices of NAICS value added for 55 industries that together comprise GDP. We focus on the 1984 to 2007 period so that our sample spans the so-called "Great Moderation" time period. Thus we have 26 time series observations. The data is logged and first differenced to ensure the series are stationary. The maximum number of factors permitted is 8, and as before we use

Table 3: Unbounded factor loadings (Case 2) with $r = 1$ (percentage of selection).

| | | Without standardization | | | | With standardization | | | | Minimum rule | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T$ | $N$ | k=0 | k=1 | k=2 | k=3 | k=0 | k=1 | k=2 | k=3 | k=0 | k=1 | k=2 | k=3 |
| 25 | 25 | 0 | 100 | 0 | 0 | 0 | 58 | 42 | 0 | 0 | 100 | 0 | 0 |
| 50 | 25 | 0 | 100 | 0 | 0 | 0 | 17 | 83 | 0 | 0 | 100 | 0 | 0 |
| 100 | 25 | 0 | 100 | 0 | 0 | 0 | 1 | 99 | 0 | 0 | 100 | 0 | 0 |
| 25 | 50 | 0 | 100 | 0 | 0 | 0 | 60 | 40 | 0 | 0 | 100 | 0 | 0 |
| 50 | 50 | 0 | 100 | 0 | 0 | 0 | 45 | 55 | 0 | 0 | 100 | 0 | 0 |
| 100 | 50 | 0 | 100 | 0 | 0 | 0 | 2 | 98 | 0 | 0 | 100 | 0 | 0 |
| 25 | 100 | 0 | 100 | 0 | 0 | 0 | 75 | 25 | 0 | 0 | 100 | 0 | 0 |
| 50 | 100 | 0 | 100 | 0 | 0 | 0 | 45 | 55 | 0 | 0 | 100 | 0 | 0 |
| 100 | 100 | 0 | 100 | 0 | 0 | 0 | 7 | 93 | 0 | 0 | 100 | 0 | 0 |

the $\text{IC}_{p2}(k)$ criterion. Evidently, standardization increases the estimated factor number in all sub-samples considered.

Table 4: Estimated Factor Number to Industry Value Added using $\text{IC}_{p2}(k)$

| Sample: | 1984–2007 | 1985–2007 | 1984–2006 |
|---|---|---|---|
| No standardization | | | |
| 1 | 1 | 1 | |
| Standardization | | | |
| 4 | 5 | 3 | |

## 5. CONCLUSION

Excessive heteroskedasticity in the finite sample can hamper factor model decompositions (such as PC) as well as associated factor-dimension selection criteria. Standardization is a common treatment for this problem. In this paper, we find that heteroskedasticity due to large factor loadings does not cause inconsistency in PC estimation as long as the signal is sufficient in the common component. In contrast, heteroskedasticity in the idiosyncratic errors may cause over-estimation of the factor number. We demonstrate that standardization can solve the over-estimation problem due to idiosyncratic heteroskedasticity, but it

can create a new source of inconsistency if the factor loadings are the source of the heteroskedasticity. This means that standardization does not necessarily result in consistent estimation, especially because the components are not separately observed so we cannot tell what is the source of heteroskedasticity. We suggest a minimum rule that can provide a simple method to make a better choice between whether or not to standardize the panel.

Although our examples and theorems specifically focus on the IC criteria and the standard principal component estimator, the proofs rely on expressions for the eigenvalues of the variance covariance matrix of the treated panels. Hence the main results are likely to be equally applicable to other eigenvalue decomposition estimation and model selection methods.

## REFERENCES

Bai, J. (2003). Inferential theory for factor models of large dimensions, Econometrica 71, 135–172.

Bai, J., and S. Ng (2002a). Determining the number of factors in approximate factor models, Econometrica 70, 191–221.

Bai, J., and S. Ng (2002b). Determining the number of factors in approximate factor models, Errata.
http://www.columbia.edu/~sn2294/papers/correctionEcta2.pdf.

Bai, J., and S. Ng (2006). Evaluating latent and observed factors in macroeconomics and finance, Journal of Econometrics 113:1-2, 507–537.

Bai, J., and S. Ng (2008). Forecasting economic time series using targeted predictors, Journal of Econometrics 146, 304–317.

Bajari, P. and L. Benkard (2005). Hedonic price indexes with unobserved product characteristics, Journal of Business and Economic Statistics 23, 61–75.

Bernanke, B. and J. Boivin (2003). Monetary policy in a data rich environment, Journal of Monetary Economics 50, 525–546.

Bernanke, B., J. Boivin, and P. S. Eliasz (2005). Measuring the effects of monetary policy: A factor-augmented vector autoregressive (FAVAR) approach, The Quarterly Journal of Economics 120, 387–422.

Boivin, J., and S. Ng (2006). Are more data always better for factor analysis, Journal of Econometrics 132, 169–194.

Chamberlain, G., and M. Rothschild (1983). Arbitrage, factor structure, and mean variance analysis on large asset markets, Econometrica 51, 1281–1304.

Greenaway-McGrevy, R., C. Han, and D. Sul (2012a). Estimating the number of common factors in serially dependent approximate factor models, mimeo, Korea University.

Greenaway-McGrevy, R., C., Han, and D. Sul (2012b). Estimation of world interest rate: The role of standardization in the estimation of common factors, mimeo, University of Texas at Dallas.

Kapetanios, G, and M. H. Pesaran (2007). Alternative approaches to estimation and inference in large multifactor panels: Small sample results with an application to modelling of asset returns, in G. Phillips and E. Tzavalis (eds.), The Refinement of Econometric Estimation and Test Procedures: Finite Sample and Asymptotic Analysis, Cambridge University Press, Cambridge.

Oulton, N. (2008). Chain indices of the cost-of-living and the path-dependence problem: An empirical solution, Journal of Econometrics 144, 306–324.

Stock, J. H., and M. W. Watson (2002). Macroeconomic forecasting using diffusion indices, Journal of Business and Economic Statistics 20, 147–62.

Yin, Y. Q., Z. D. Bai, and P. R. Krishnaiah (1988). On the limit of the largest eigenvalue of large dimensional sample covariance matrix, Probability Theory and Related Fields 78, 509–521.