

FACTOR RESIDUALS IN SUR REGRESSIONS: ESTIMATING PANELS ALLOWING FOR CROSS SECTIONAL CORRELATION

DONALD ROBERTSON AND JAMES SYMONS

ABSTRACT. This paper describes a method for estimating panels by imposing a factor structure on the residuals. The method allows SUR estimation of panel models by providing a full-rank estimator of the system covariance matrix when the usual estimate is rank-deficient. We characterise completely the circumstances when this is possible. When the usual estimator is of full rank, our procedure provides a more parsimonious representation of the covariance matrix, which can lead to efficiency gains in finite samples. Monte Carlo analysis of convergence regressions and PPP regressions in the Heston-Summers data-set indicates that the proposed estimator has better performance in terms of RMSE and bias than standard panel or SUR estimators (where available), as well as offering unbiased inference.

1. INTRODUCTION

The increasing availability of data structures with both time-series and cross-sectional dimensions and the possibility of overcoming the difficulties associated with both pure time series and pure cross sectional work have given renewed interest to the econometric issues that arise with such structures. In this paper we propose an alternative method for the estimation of regression models on such data-sets.

The standard model characterises the relationship between the dependent variable and explanatory variables as a linear regression with random errors

$$(1) \quad y_{it} = \alpha_i + \beta' x_{it} + u_{it} \quad i = 1, \dots, n, \quad t = 1, \dots, T$$

Depending on the data structure available, further assumptions are generally made to identify the parameters in (1). If T is sufficiently large and the assumption made that the errors are uncorrelated over time (i.e. $E(uu') = \Sigma$ where $u' = (u_{1t}, u_{2t}, \dots, u_{nt})$ for all t), the system can be treated as seemingly unrelated regression equations (Zellner, 1962) and α_i , β , and the variance-covariance matrix of the errors Σ can be straightforwardly estimated. If, further, the errors are assumed to be jointly normally distributed, estimates are available by maximum likelihood.

If T is not sufficiently large, and in particular if $n > T$, then such an approach is not feasible, because the usual estimator of the error covariance matrix is rank-deficient. There are a variety of techniques discussed in the literature for dealing with the $n > T$ case. One method is to assume that Σ is diagonal. If it is further

Date: November 25th 2000.

Key words and phrases. Panel data. Cross sectional correlation. Factor analysis.

Part of this work was completed whilst Robertson was a Jean Monnet Fellow at the European University Institute, Florence. We are grateful to Hashem Pesaran, Nick Rau and Ron Smith for helpful comments and suggestions, and to seminar participants at Bristol, Sussex and EUI. All remaining errors are our own.

assumed that Σ is a scalar matrix, one has the usual panel estimator. This restricted structure is the one-way error components model. The parameters α_i can be treated as fixed or random, and the model estimated by OLS, GLS, or, if one is willing to specify a joint distribution for the u_{it} and the α_i , maximum likelihood. In these circumstances one can also allow for time fixed effects (the two-way error components model). If the explanatory variables include a lagged dependent variable, estimation is complicated in the fixed effects formulation: see Nickell (1981). In both SUR and panel estimation, it is possible to allow for some restricted serial correlation in the error process (see for example Liang and Zeger (1986) for the static case, Keane and Runkle (1992) for the dynamic case). If one relaxes the assumption that the error variances are equal, one can weight observations accordingly to obtain WLS estimators as in Fisher (1993). For Σ non-diagonal, one could always collapse some of the cross-sectional units into group averages to make $n < T$, as is often done for ML estimation of the factor model in finance. A further possibility is to replace the required inverse of Σ by its generalised inverse. Under certain conditions this is a BLUE estimator (see Theil, 1971), and brings the $n > T$ case within the range of SUR estimation, though we know of no applications of this technique in a panel context. The advantage of the factor approach we propose *vis-à-vis* this possibility lies in its parsimony: we shall see that the factor approach is much more efficient in small samples, even in the full-rank case, at least for the cases we consider¹.

The assumptions on the correlation structure Σ are usually made for identifiability reasons, rather than descriptive accuracy. In particular the panel approach imposes zero correlation on the shocks to the cross-sectional units at each date so that estimation can proceed. The hope, presumably, is that by conditioning on sufficiently many explanatory variables x_{it} , what remains can be treated as a purely idiosyncratic shock, uncorrelated over time and the cross-sectional dimension, though with perhaps a t specific effect and a i specific effect. In practice, if the individual units are countries, firms or even individuals, this is unlikely to be the case. The econometrician rarely has sufficient explanatory variables to remove all correlated behaviour, and what is left over may not be well described by an $\alpha_i + f_t + \epsilon_{it}$ formulation². Ideally one would like to allow for some sort of correlation pattern across the shocks. In a time series, it is natural to specify the correlations between disturbances to be functions of distance, measured by time, and a small set of parameters (such as in an ARMA process). In a cross-section, however, there is often no unambiguous concept of distance. One approach is to define a metric on the cross-section using some notion of physical distance and allow for a correlation structure with this measure playing the part of separation in a time series (see e.g. Conley (1999) who discusses the consistency of GMM estimation in these circumstances). Thus, for example, with a data-set consisting of many countries, one could allow for correlated shocks on adjacent countries, with zero correlation at greater distances. This approach is predicated on formulating an appropriate concept of nearness but clearly one country may be near a second for some purposes, but not for others. In the SUR framework the possibility of correlated shocks is allowed for through the covariance matrix Σ . However estimation of this matrix can often be

¹It is to be emphasised that this is a small sample property. It is well known that, with a minor caveat, imposing restrictions on Σ brings no efficiency gains in large samples. See Greene, 1990.

²Frees (1995) discusses testing cross sectional correlation in panel data

expensive in terms of degrees of freedom, and a poor estimate of this matrix may contaminate the estimation of the parameters of interest in 1.

In this paper we propose a method that retains the flexibility of the SUR approach in allowing for correlated shocks, yet is more parsimonious than SUR, so can be expected to be more efficient in general. It can be implemented when $n > T$, and provides a more general specification than the panel approach. We impose a factor structure on the covariance matrix and estimate the factors by maximum likelihood techniques. This procedure can give an improvement in both bias and RMSE of estimators. The estimation procedure is essentially feasible GLS and so has the usual efficiency properties for large T .³ In Section 2 we set out the details of the method. Section 3 applies the method to Barro-type convergence regressions and PPP regressions in the Heston-Summers data-set. Sections 4 and 5 assesses the performance of the factor approach by Monte Carlo and suggest some approximate methods. Section 6 concludes.

2. A FACTOR RESIDUAL APPROACH

We specify the model as follows

$$(2) \quad y_t = X_t\beta + u_t \quad t = 1, \dots, T$$

where y_t and u_t are $n \times 1$ vectors, X_t is an $n \times v$ matrix of explanatory variables observed at t , β is a vector of unknown parameters to be estimated, u_t is a vector white noise process with $E(u_t u_t') = \Sigma$, and $E(X_{ijt} u_{kt}) = 0$ all i, j, k, t .

We make the further assumption that

$$(3) \quad \Sigma = \Lambda\Lambda' + \Psi$$

where Λ is a $n \times m$ matrix of so-called factor loadings and Ψ is a diagonal $n \times n$ matrix with diagonal elements $\psi_1, \psi_2, \dots, \psi_n$, where $\psi_i > 0$ reflects idiosyncratic effects. This allows for some contemporaneous correlation between shocks, expressed as a function of fewer parameters than the unconstrained Σ if $m < n$. Note that the factor model generalises fixed effects models directly, as the fixed effects can be entered as elements of X . A time random-effects model is equivalent to a one-factor model with Ψ diagonal and Λ proportional to $(1, 1, 1, \dots, 1)'$. Ψ scalar and $m = 0$ gives the usual panel formulation.

The factor model has an attractive interpretation because it amounts to specifying that the residuals take the form $u_t = \Lambda_1\phi_t^1 + \Lambda_2\phi_t^2 + \dots + \Lambda_m\phi_t^m + \varepsilon_t$ where $E(\phi_t^i \phi_s^j) = \delta_{ij}\delta_{ts}$ and $E(\varepsilon_t \varepsilon_s') = \delta_{ts}\Psi$ and Λ_k is a column vector of weights. The ϕ_s can be interpreted as m common shocks and the elements of each Λ_k give the loading or impact of each of these shocks on the cross sectional units. For example, if the units are economies, the first factor ϕ_t^1 might represent a world demand shock with Λ_1 as the relative openness of the economies, the second factor ϕ_t^2 an agricultural shock etc. For $m = n$ such a decomposition can trivially be obtained; the usefulness of this decomposition arises when m may be taken to be much smaller than n . These common shocks provide cross-sectional correlation in the error structure, with the ε_t adding an idiosyncratic term.

Given an estimate S of Σ , the likelihood function up to an additive constant is

$$(4) \quad L(\Sigma) = L(\Lambda, \Psi, S) = -\frac{T}{2} (\log \det \Sigma + tr(\Sigma^{-1}S))$$

³For fixed T , the estimators as n grows may or may not be consistent.

(see Lawley and Maxwell, 1963).

The first-order conditions are (Lawley and Maxwell or Jöreskog, 1967)

$$(5) \quad \begin{aligned} \frac{\partial L}{\partial \Lambda} &= -2\Sigma^{-1}(\Sigma - S)\Sigma^{-1}\Lambda = 0 \\ \frac{\partial L}{\partial \Psi} &= -\text{diag}(\Sigma^{-1}(\Sigma - S)\Sigma^{-1}) = 0 \end{aligned}$$

Classical factor analysis assumes that S is of full rank, and thereby obtains estimates of Λ and Ψ by maximum likelihood. The results cannot immediately be generalised to the rank-deficient case. There are methods in the literature for dealing with this situation; for instance, Maxwell (1981) suggests estimation via norm-minimisation⁴ It is also possible to obtain factors as principal components, even in the rank-deficient case. In the finance literature, Connor and Korajczyk (1986, 1988) estimate factors using the Chamberlain and Rothschild (1983) method of asymptotic principal components in the rank-deficient case. These factors are identified for asymptotically large n given certain assumptions. The advantage of maximum likelihood is that one has, in large samples at least, a natural basis for inference and model selection. The technical contribution of this paper is to show that, subject to certain restrictions on the number of fitted factors, the likelihood function is bounded and attains its bound at an invertible Σ , irrespective of the rank of S . Hence maximum-likelihood estimators of Λ and Ψ can be obtained from direct maximisation of the likelihood function, even in the rank-deficient case⁵. The likelihood is invariant to multiplication of Λ by a unitary matrix so the matrices Λ and Ψ are identified only up to such a matrix. This can be resolved by fixing Λ to have zeros above the diagonal, thus removing the difficulty in *estimation*; though if one wished to offer *interpretation* of the factors it would, of course, be problematical.

We summarise our results concerning the likelihood function in the following. If S has full rank there is no problem in maximising the likelihood (see e.g. Jöreskog, (1967)). We show in the appendix that, as long as one does not attempt to fit too many factors, the likelihood function for the general factor model is bounded and attains its bound, even if the initial estimate of the variance-covariance matrix is rank-deficient. The estimator of the covariance matrix thus obtained has full rank. This result clearly has applications beyond those we explore below. To obtain estimates of the factors we can therefore operate directly on the likelihood function (4). As is well known in the case S has full rank, it can be difficult to find the maximum of the likelihood function numerically. The situation is the same when S is rank-deficient. The geometry of the likelihood function is explored further in the appendix.

In the situation where the likelihood is bounded, the general shape is akin to the positive orthant of a multi-dimensional volcano. The maximum is obtained on the rim or on the boundary of the admissible region (these are known as Heywood solutions). Clearly search algorithms will usually have great difficulty with this

⁴The norm in question being $\|B\| = (\text{tr}B'B)^{\frac{1}{2}}$, so that one would minimise

$$\|S - \Sigma\|^2 = \sum_{i,j} (S_{ij} - \Sigma_{ij})^2.$$

⁵Given the factor loading matrix Λ and Ψ , the factors can then be recovered by regression, if desired.

surface. Motivated by this, we propose a mixture of search algorithms to ensure convergence. The numerical technique is discussed in detail in the appendix.

Having obtained estimates of the factors, estimates of β can then be obtained by feasible GLS, by minimising $\sum_{t=1}^T (y_t - X_t\beta)' \hat{\Sigma}^{-1} (y_t - X_t\beta)$, where $\hat{\Sigma} = \hat{\Lambda}\hat{\Lambda}' + \hat{\Psi}$ to give $\hat{\beta} = \left(\sum_{t=1}^T X_t' \hat{\Sigma}^{-1} X_t \right)^{-1} \sum_{t=1}^T X_t' \hat{\Sigma}^{-1} y_t$. One could in principle iterate this approach (as is usually done in a SUR framework) to obtain full maximum likelihood estimates. Standard errors can be obtained from the usual GLS formula.

Asymptotic consistency results are complicated by the incidental parameter problem (Neyman and Scott, 1946) arising particularly in this context when time or individual fixed effects are included in the model. For T -asymptotics with fixed n , if we begin with an initial T -consistent estimate of Σ , the estimates of Λ and Ψ will be T -consistent if the factor model is true⁶. Hence our GLS estimator of β inherits the T -asymptotic properties of GLS. Thus, whereas estimates of time fixed effects will not be T -consistent, other parameter estimators typically will be consistent. The gain here is a more parsimonious representation of the covariance matrix, leading, as we shall see, to potential efficiency gains even when T is moderate.

The case of n -asymptotics is a little more murky. Clearly one could not expect to obtain n -consistent estimates of individual fixed effects, for given T . Beyond this, the Nickell result shows that, in the presence of individual fixed-effects, estimates of the structural parameters can be n -inconsistent. If $\dim(\beta)$ is fixed as n grows, however, one might hope to say something, but we have no results as yet. One could as well seek asymptotics for large T but larger n . It is true that, if the data are fixed in repeated samples and S is obtained by OLS on each of the n equations, then the GLS estimate of $\hat{\beta} - \beta$ will be an odd function of the residuals, and hence median-unbiased. However, irrespective of these considerations, we shall see below that, for fixed sample size, the proposed estimator has better performance measured by standard loss functions than both conventional SUR (where feasible) and panel methods, and in this sense can provide superior estimates. The empirical contribution of this paper is thus to suggest the use of this factor-residual model in a SUR framework, permitting estimation of panel models allowing for cross-sectional dependence. We illustrate with an empirical example.

3. AN APPLICATION: CONVERGENCE REGRESSIONS

We illustrate the above techniques by applying the method to study convergence in the Heston-Summers data-set. We seek to measure convergence in GDP per capita relative to the US. This problem has been studied by many authors (Barro, 1991, Mankiw Romer Weil, 1992, Evans and Karras, 1996, Islam, 1995, Lee, Pesaran and Smith, 1995, Caselli *et al.*, 1996, among others) mainly using the Heston-Summers database. Barro-type regressions correlate average growth rates and initial values; others attempt to control for extra right-hand-side variables such as human capital, savings rates etc. and study conditional convergence. A typical Barro result is convergence at about 2% per annum; more recently a number of panel studies have suggested a somewhat higher rate of convergence (see Evans and Karras, 1996, Islam, 1995, Lee, Pesaran and Smith, 1995, Caselli *et al.*,

⁶By application of Slutsky's Theorem, given that, as is easily shown, the mapping $S \rightarrow \hat{\Sigma}_{Factors}$ is continuous and takes the value Σ at $S = \Sigma$.

1996). We exploit the full time-series and cross-sectional dimensions of the Heston-Summers data-set, and study unconditional convergence, though our methodology could straightforwardly be applied to study conditional convergence. Imposing a factor structure means we can allow for the shocks that disturb the system to be correlated across countries – for instance oil price shocks or financial crises. We specify the model as

$$(6) \quad y_{it} = \rho y_{it-1} + \varepsilon_{it}$$

where $y_{it} = \ln(\text{GDP per capita as percentage of US})$. We shall impose a factor structure on the ε_{it} , allowing for correlated shocks. Note that we do not allow for individual fixed-effects. As discussed above, these would cause problems in a dynamic regression if T were small (see Nickell, 1981), though there is a developing literature addressing this problem by applying GMM techniques to the differenced equation (see Islam, and Caselli *et al.*). In an unconditional convergence regression, the natural formulation does not involve individual intercepts.⁷ In estimation we include a simple intercept.

3.1. Results and selection of number of factors. We estimate this model for the subset of the OECD countries ($n = 22, T = 41$) and for the whole data-set ($n = 103, T = 31$). We also estimate (6) using a panel estimator (i.e. applying OLS to the stacked data, appropriate when Σ is a scalar matrix) for both the OECD and world subsets, and using SUR for the OECD countries. One key issue is the selection of the number of factors to be fitted. Likelihood ratio tests are available, as well as information criteria such as Akaike, Schwarz-Bayes or Hannan-Quinn⁸. We discuss the performance of these information criteria later.

For the OECD subset, the estimate S of the matrix Σ indicates considerable correlated structure in the errors. The average modulus of the off-diagonal correlations in this matrix is 0.375. With correlations of this magnitude, panel/OLS estimation will produce biased standard errors, and be less efficient than unrestricted SUR (see Di Liberto & Symons, 1998) The estimation results are set out in Table 1.

⁷Though in principle there would be no difficulty in combining any fixed effects estimator with our factor residual model by premultiplying (6) by $\widehat{\Sigma}^{-1/2}$, as long as an initial consistent estimate of Σ is available.

⁸These criterion function are of the form $\text{LogLikelihood} - f(T) \cdot (\#params)$ where $f(T) = 1$ for Akaike, $f(T) = \frac{1}{2} \ln(T)$ for Schwarz, and we use $f(T) = \ln(\ln(T))$ for Hannan-Quinn.

Fitted Σ	intercept	$\hat{\rho}$	log likelihood	AIC	SBC	HQC
scalar(OLS)	0.12526 (.01314)	0.97234 (.00323)	0	-1.00	-1.84	-1.31
0-factor(WLS)	0.12847 (.01321)	0.97167 (.00322)	60.01	38.01	19.43	31.29
1-factor	0.11582 (.01191)	0.97422 (.00277)	266.36	222.36	185.20	208.92
2-factor	0.11702 (.01125)	0.97463 (.00260)	305.30	240.30	185.41	220.45
3-factor	0.11445 (.01090)	0.97527 (.00252)	322.34	237.34	165.56	211.38
4-factor	0.12908 (.01191)	0.97175 (.00278)	339.37	235.37	147.55	203.62
5-factor	0.13297 (.01198)	0.97085 (.00279)	354.44	232.44	129.42	195.19
6-factor	0.13316 (.01181)	0.97071 (.00279)	370.93	231.93	114.56	189.50
7-factor	0.12902 (.01178)	0.97169 (.00280)	382.77	227.77	96.88	180.44
8-factor	0.12896 (.01109)	0.97176 (.00264)	394.44	224.44	80.89	172.54
19-factor	0.13239 (.00924)	0.97076 (.00221)	445.77	176.77	-50.39	94.64
20-factor	0.13239 (.00924)	0.97076 (.00221)	445.77	173.77	-55.92	90.72
n -factor(SUR)	0.13239 (.00924)	0.97076 (.00221)	445.77	192.77	-20.88	115.52

Table 1: Factor Residual, SUR and OLS/Panel Convergence Regressions for the OECD subset of Heston-Summers

Notes: (i) Σ scalar is Panel/OLS; 0-factor is weighted OLS; n -factor is SUR. (ii) Log-likelihood is normalised to be 0 for Σ =scalar (OLS). (iii) Since Λ is restricted to be upper diagonal, adding a further factor adds fewer than 22 parameters. (iv) Nominal standard errors in parentheses, calculated as $\hat{\sigma}^2(X'\hat{\Sigma}^{-1}X)^{-1}$ for $\hat{\Sigma}$ in each row. These standard errors are fairly dubious *a priori* for a process so close to the unit circle.

The results are listed by decreasing degree of restriction on the covariance matrix. If Σ is restricted to be a scalar matrix, one has OLS; if Σ is diagonal (0-factors) one has weighted least squares; each factor model is nested in a higher-order model; SUR (n -factor) is the least restricted. OLS and WLS are strongly rejected against models with non-zero covariance structure by likelihood ratio tests. The restriction from seven to six factors would not be rejected at the 5% level by likelihood ratio, but further restrictions would be rejected. Using the information criteria, Akaike, Schwarz and Hannan-Quinn all select the two-factor model. Likelihood-ratio tests of constant size will not lead to consistent order selection, unlike SBC and HQC.

For the entire Heston-Summers data-set ($n = 103, T = 31$) we obtain the results in Table 2.

Fitted Σ	intercept	$\hat{\rho}$	log likelihood	AIC	SBC	HQC
scalar(OLS)	-0.01774 (.00398)	1.00617 (.00134)	0	-1.00	-1.70	-1.22
0-factor(WLS)	-0.01581 (.00321)	1.00564 (.00091)	626.91	523.91	451.75	500.83
1-factor	-0.01953 (.00326)	1.00608 (.00076)	1001.15	795.15	650.83	748.98
2-factor	-0.01456 (.00310)	1.00462 (.00073)	1145.27	837.27	621.49	768.24
3-factor	-0.01684 (.00300)	1.00589 (.00073)	1274.88	865.88	579.34	774.21
4-factor	-0.01720 (.00286)	1.00597 (.00071)	1376.93	867.93	511.32	753.85
5-factor	-0.01681 (.00280)	1.00588 (.00069)	1474.56	866.56	440.60	730.29
6-factor	-0.01692 (.00260)	1.00595 (.00066)	1573.85	867.85	373.23	709.62

Table 2: Factor Residual and OLS/Panel Convergence Regressions for the Heston-Summers Data-set.

Notes: As for Table 1, except that SUR is omitted since $n > T$.

Akaike is fairly flat between the three, four and five factor models for the world regressions, Hannan-Quinn selects three factors and Schwarz points to a one-factor model, reflecting the heavier weight placed on extra parameters by this criterion. Again a scalar or diagonal covariance matrix is strongly rejected.

The parameter estimates are in line with those obtained by other authors, and the literature on convergence clubs. We find a point estimate of the rate of convergence of approximately 2.5% per annum for the OECD countries, and no convergence for the full Heston Summers data-set.

3.2. Confidence Intervals for the Convergence Parameters. We construct confidence intervals for the three estimates of the convergence parameter for the OECD subset by Monte Carlo. We define a 95% confidence interval as the points lying between the 2.5 percentile and 97.5 percentile points of the sampling distribution obtained from 1000 repetitions of each estimation procedure applied to artificial data generated by the relevant estimated parameters in Table 1 and the *original* estimated variance-covariance matrix S to generate the errors⁹. For the factor model, we select the estimated parameter values from the (preferred) two-factor model for the Monte Carlo, and always estimate two-factor models on the generated data. The results are set out in Table 3.

⁹That is, we do not generate the data using an underlying factor residual model. This will ensure that we are not biasing the Monte Carlo in favour of such models.

	true parameter	2.5 percentile	mean	97.5 percentile
OLS				
α	0.12526	0.09350	0.13105	0.17306
ρ	0.97234	0.96097	0.97093	0.98000
SUR				
α	0.13239	0.10678	0.13614	0.15905
ρ	0.97076	0.96198	0.96986	0.97686
Factors				
α	0.11702	0.09619	0.11755	0.14025
ρ	0.97463	0.96947	0.97449	0.97927

Table 3: Empirical Confidence Intervals from 1000 repetitions.¹⁰

It will be observed that the 95% confidence intervals for the factor model are about half the width of the OLS confidence intervals and two-thirds of the SUR value. It is useful to compare these intervals with the apparent confidence intervals (estimate plus or minus two estimated standard errors) from Table 1. Apparent confidence intervals for OLS are expected to be wrong, even asymptotically, when Σ is non-diagonal, while SUR will be correct in large samples but wrong in small samples when an estimate of Σ is used in place of its true value to compute standard errors. One finds that the apparent confidence interval width by OLS is about 70% of its true value, for SUR about 60%. Factors, in contrast, estimates the true confidence interval almost exactly. For this example at least, the factor method offers much superior inference. We return to this issue below.

3.3. A Further Application: PPP. We illustrate the technique with a further application, looking at PPP in the Heston Summers data-set. The existing literature finds mixed evidence for PPP using time-series, cointegration and panel methods. Using real exchange rate data from Heston Summers, we have $n = 103, T = 31$ (corresponding to the years 1960-1990). We take as our dependent variable $\ln(\text{exchange rate})$ and normalise by subtracting from each time series its average value over 1964-69. The intention is thus to model deviations of exchange rates from some “normal” level, here taken as the average over the second half of the sixties. This is a compromise between allowing for individual country fixed effects (equivalent to subtracting time averages of the data for each country), entailing problems of bias in a dynamic panel, and normalising the exchange rates on a particular year which has the risk of choosing an anomalous reference point. We again estimate an AR(1) regression, fitting an intercept (which should be zero if PPP holds). The average modulus of the correlations in these data is 0.275, indicating considerable cross-sectional dependence. Thus, as for the convergence regressions, conventional panel methods are likely to lead to distorted inference and loss of efficiency.

The results are in Table 4.

¹⁰In fact, SUR and OLS are based on 5000 repetitions, factors 1000.

Fitted Σ	intercept	$\hat{\rho}$	log likelihood	AIC	SBC	HQC
scalar(OLS)	0.00205 (.00249)	0.90534 (.00911)	0	-1.00	-1.70	-1.22
0-factor(WLS)	-0.00015 (.00161)	0.95222 (.00776)	805.43	702.43	630.27	679.34
1-factor	-0.01809 (.00197)	0.95810 (.00722)	1596.06	1390.06	1245.73	1343.89
2-factor	0.00767 (.00227)	0.86290 (.00930)	1835.08	1527.08	1311.29	1458.05
3-factor	0.00344 (.00237)	0.85584 (.00945)	2000.75	1591.75	1305.21	1500.08
4-factor	0.00328 (.00228)	0.85887 (.00929)	2054.18	1545.18	1188.58	1431.10
5-factor	0.00262 (.00224)	0.86466 (.00901)	2173.98	1565.98	1140.02	1429.71
6-factor	0.00191 (.00217)	0.84166 (.00975)	2272.95	1566.94	1072.32	1408.71

Table 4: Factor Residual and OLS/Panel PPP Regressions for the Heston Summers Dataset.

Note: Nominal standard errors in parentheses.

Both Akaike and Hannan-Quinn select a three factor model, whereas Schwarz marginally prefers two factors. Whilst panel/OLS gives an AR parameter of 0.91, allowing for the correlated error structure reveals $\hat{\rho} = 0.86$. The difference between the half-lives of deviations from PPP implied by these estimates is of economic significance.

4. MONTE CARLO ANALYSIS OF THE METHOD.

4.1. Deciding on the number of factors. We take the fitted three-factor covariance matrix from the Heston Summers PPP regressions and use this to generate 103×31 matrices of correlated normal variates. These are used to generate artificial AR(1) time series of length 31, using the parameters in Table 1. An AR(1) model is then fitted to each of the 103 time series and the residuals used to form a rank-deficient 103×103 variance-covariance matrix. AIC, SBC & HQC are then calculated to see whether they are able to identify the correct number of factors, here three. This procedure is rather time-intensive: the five- and especially the six-factor models are automatically heavily overparameterised and this seems to cause convergence to be very slow. We therefore limit the number of repetitions to 100. Table 5 gives the frequency of selection of different numbers of factors by the three information criteria.

Frequency of fitted factors:	1	2	3	4	5	6
Chosen by AIC	0	1	95	4	0	0
Chosen by SBC	0	19	81	0	0	0
Chosen by HQC	0	1	98	1	0	0

Table 5. Performance of various information criteria in identifying correct factor model.

Results from 100 repetitions of underlying 3-factor model.

The HQC does astonishingly well, getting the number of factors correct 98/100. AIC is nearly as good. SBC tends often to choose smaller models.

4.2. Factors vs SUR vs OLS - who does best? We generate Monte Carlo replications of the data using as initial values the 1950 values of log GDP for the OECD subset of Heston-Summers. We generate series for y_{it} via

$$y_{it} = \alpha + \rho y_{it-1} + \varepsilon_{it}$$

where α and ρ are the estimated parameters from the SUR of this equation and the ε_{it} are correlated normally distributed random numbers with variance-covariance matrix as estimated from the residuals of the SUR. We use an empirical variance-covariance matrix for these Monte Carlo, rather than generating data with a given factor structure, to ensure that the experiments are not biased in favour of factors. We generate time series of different lengths T holding $n = 22$. On these generated data we estimate AR(1) convergence regressions by OLS/Panel, SUR (whenever $n < T$) and the factor approach, assuming a two-factor model, as suggested by Table 1, a one-factor model and a zero-factor model (WLS). Given that time fixed- or random-effects will generate a model close to the one-factor specification, it is of interest to check the performance of the factor specification against a model where cross-sectional dependence is modelled as time fixed-effects and this is also calculated. The root-mean-squared-error and bias from 1000 repetitions are graphed in Figures 1 and 2. Note first that there is nothing between the one- and two-factor models. The factor-residual approach beats OLS/Panel and SUR in RMSE and bias for $n = 22$ and T running from 5 to about 60. One- and two-factor models outperform time fixed-effects in RMSE, with similar bias performance. Eventually, as T rises, SUR starts to dominate, though of course this is holding the number of factors fixed. If the number of factors were allowed to rise as T grew, this may no longer hold. Since in these Monte Carlo experiments the variance-covariance matrix generating the data is, in fact, a genuine empirical matrix and is thus presumably a 22-factor model¹¹, a consistent selection criterion such as SB or HQ would ultimately select $m = 22$ and fit the model as well as SUR. Thus the superior performance of factors for T up to 60 illustrates the efficiency gain of parsimony, even, as here, at the expense of misspecification.

Whilst one- and two-factor models perform about equally well in RMSE terms, rather unexpectedly, we find in Figure 1 that WLS does worse even than OLS for some sample sizes. In general, transforming the data so that the residuals are homoscedastic does not always lead to an increase in efficiency when Σ is non-diagonal.¹² This is a useful experiment because it shows that the efficiency gains from the factor approach are due to exploiting the non-diagonal nature of Σ rather than the inequality of the terms along the main diagonal. The fact that the one-

¹¹Even though, in Table 1, treating this matrix as the sample estimate of an underlying m factor model, we selected $m = 2$.

¹²Assume 2 contains just one independent variable x_t with covariance matrix $V = E(x_t x_t')$. Then OLS is the best among *all* estimators that weight the data by a diagonal matrix (thus including WLS) if

$$diag(V\Sigma) = \lambda diag(V)$$

for some $\lambda > 0$. Examination of this condition shows that for OLS to be both best and different to WLS requires both V and Σ to be non-diagonal. This apart, the condition can certainly hold, which implies that there are regions of V, Σ space where OLS is better than WLS. These results are for the case of known Σ .

and two-factor models perform equally well means presumably that the efficiency gains from parsimony are off-set by losses incurred by the imposition of rejectable restrictions on the covariance matrix.

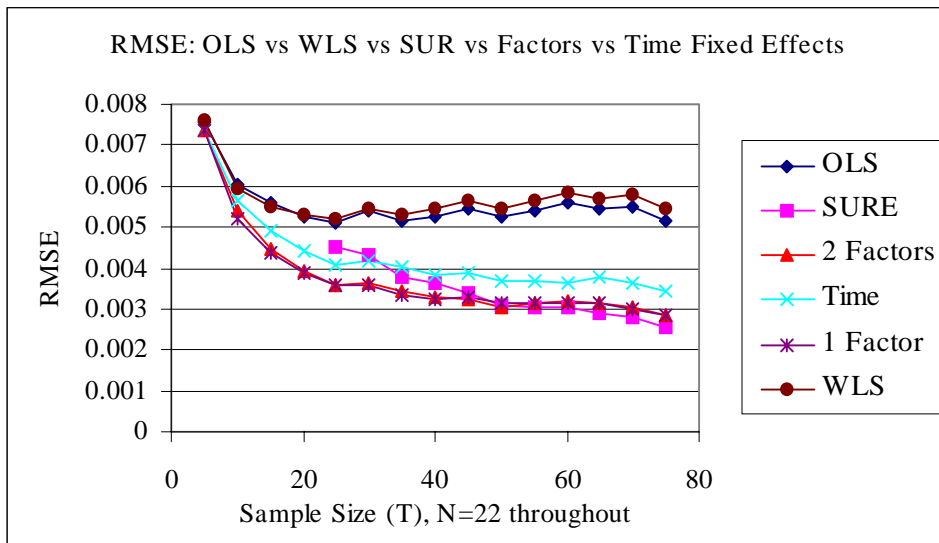


Figure 1

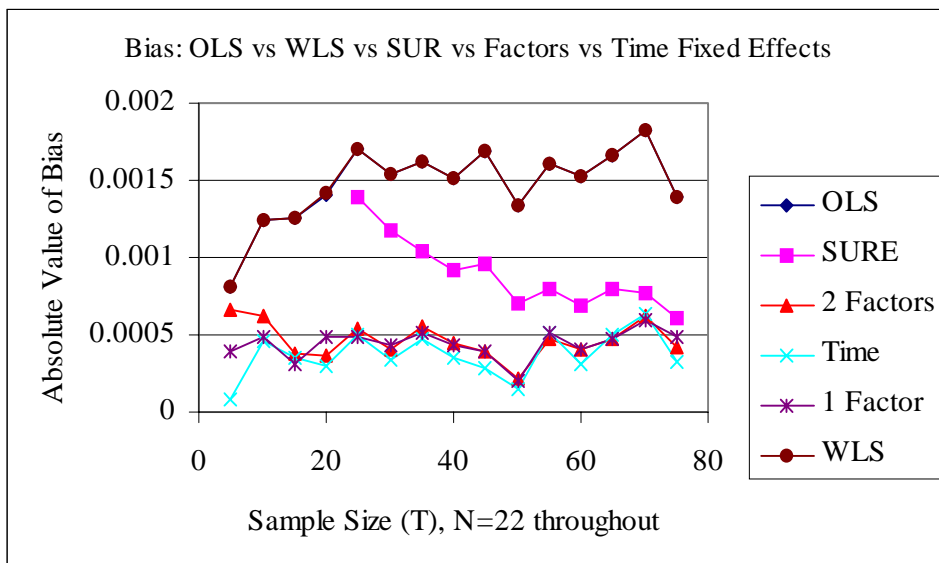


Figure 2

4.3. Actual versus calculated standard errors for the various estimation procedures. For the Monte Carlo described above we calculate the sampling standard error of the estimated parameter $\hat{\rho}$ from 1000 repetitions for the OLS, SUR and two-factor model. We also obtain the average calculated standard error from the $\hat{\sigma}^2(X'X)^{-1}$ formula for each estimator. The ratio of the estimated to the true standard error as the sample length (T) varies is graphed in Figure 3. We find that OLS and SUR provide seriously biased inference over a wide range of sample sizes, in both cases calculated standard errors substantially understating the true sampling variation of the estimator. For instance, with a sample size of 15 (so that the panel has $n=22$, $T=15$), OLS calculated standard errors are approximately half the sampling standard error from the 1000 repetitions. The factor residual approach provides a much superior guide to the true sampling variation, though still tends to underestimate somewhat. The performance of SUR when the time and cross sectional dimensions of the data are of comparable magnitude is very poor. Of course these results depend on the correlation structure in the underlying Monte Carlo, which here is drawn from the estimates of the AR(1) model for GDP per capita on the OECD subset of the Heston and Summers database. But we have no reason to think that this is a particularly exceptional database in terms of error correlation structure, so these results must cast real doubt on inference in existing panel studies¹³.

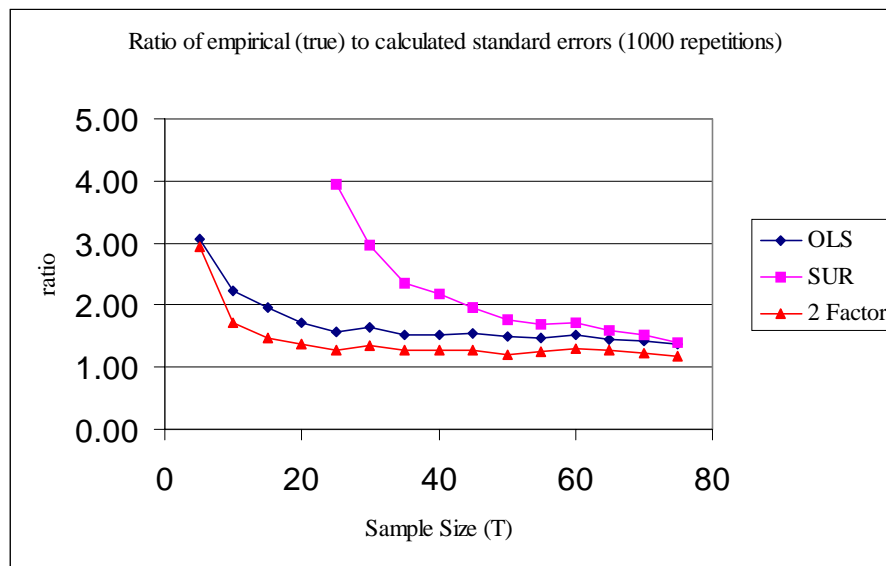


Figure 3

4.4. Normality of the distribution of the various estimators. For each of the Monte Carlo above we investigated the distribution of the 1000 estimates of the coefficient ρ . In general estimates derived from the factor residual technique showed less departure from normality than others. For example, for the sample $n = 22$, $T = 25$ we obtain the following histograms and diagnostics

¹³Driscoll and Kraay (1998) have recently given a method of obtaining asymptotically correct standard errors for OLS.

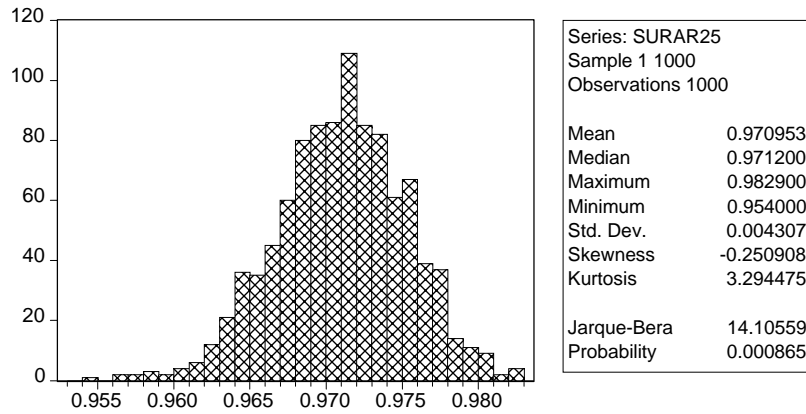


Figure 4 SUR Estimation

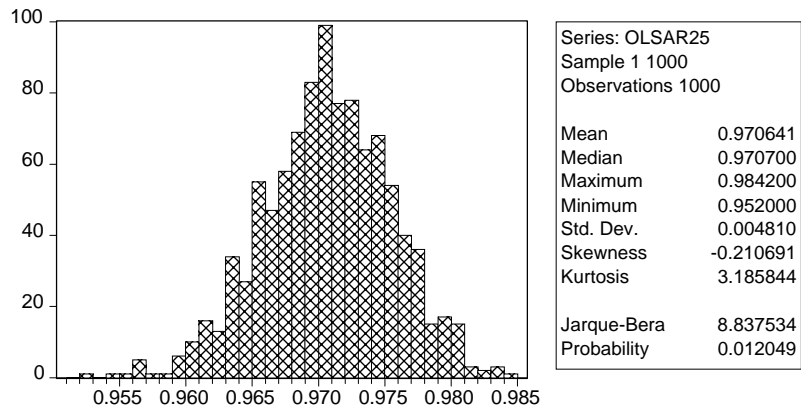


Figure 5 OLS Estimation

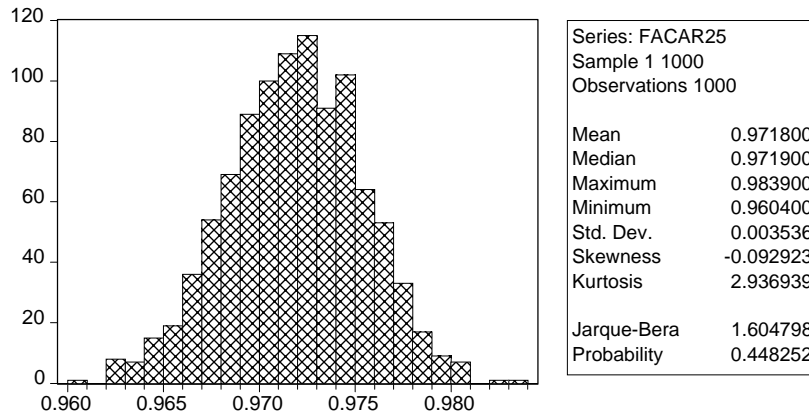


Figure 6 Factors Estimation

Notes to Figures 4, 5 and 6: Histograms of 1000 estimates of AR1 coefficient by the various estimation techniques. Sample design $n = 22$, $T = 25$ and underlying parameters of OLS estimate of AR1 convergence regression in Heston Summers OECD database.

SUR and OLS are generate a detectably non-normal distribution of estimates, whilst for the factor residual approach it is not possible to reject the null of normality of the estimator at conventional levels.

The Monte Carlo described above is for a particular set of parameter values and a particular design, and we must be wary of generalising too freely from these results. However the conclusion of these investigations is clear. The factor residual approach is superior in RMSE and bias to both OLS and SUR (where available) over a large range of sample sizes, the estimates derived from the factor residual approach follow more closely a normal distribution when the underlying errors are normal, and their estimated standard errors provide a more reliable guide for inference.

5. APPROXIMATE FACTOR TECHNIQUES

Full maximum likelihood estimation of the factor model can be rather cumbersome and time-consuming due to the nature of the likelihood surface. In this section we briefly discuss some techniques that provide an approximation to maximum likelihood. We also set out the advantages and drawbacks of these approximations as they appear to us.

5.1. Principal Components. One can easily approximate the rank deficient S matrix by its m principal components and defining the corresponding Ψ matrix to fit the diagonal exactly. This is very simple to implement, very quick and will usually produce an invertible estimator $\hat{\Sigma}$. We have experimented with such a principal component scheme, fitting the same number of principal components as the factor model selects, and find that it produces very similar results in terms of rmse etc. Thus principal components offers a straightforward method to capture correlated error structures in panel data, the difficulty being of course that, in contrast to ML, one has no obvious diagnostic for the choice of the number of principal components.

5.2. Norm Minimisation. Principal Components can be iterated to minimise $\|\Sigma - S\|$ where $\Sigma = \Lambda\Lambda' + \Psi$. For given Ψ we take the eigenvectors corresponding to the largest eigenvalues of $S - \Psi$ as our elements of Λ . We then define a new Ψ to match the diagonal elements of S exactly given this Λ . i.e. set $\Psi = \max \text{diag}(S - \Lambda\Lambda', 0)$. This procedure produces a non-increasing sequence of values of the norm and will converge to some $\hat{\Sigma}$. It is again simple to program and tends to converge quite quickly. The problem with norm-minimisation is that, unlike ML, it is not guaranteed to find an invertible $\hat{\Sigma}$. We have artificial examples of covariance matrices S where, if one factor is fitted, Heywood solutions of multiplicity¹⁴ greater than unity are found by norm-minimisation. In this case the estimated $\hat{\Sigma}$ cannot be invertible.¹⁵ This is inconvenient but not necessarily fatal as a generalised inverse could always be used for the GLS.

¹⁴The multiplicity of a Heywood solution is the number of zero components of $\hat{\Psi}$.

¹⁵An intuition for the singularity-hating behavior of ML can be obtained from consideration of (4): as an eigenvalue of Σ approaches zero, the second term on the right contributes a negative infinity of higher order than the positive infinity contributed by the first term. This tends to cause ML to choose non-singular Σ s. The norm, in contrast, has no such reason to shun singular matrices.

6. CONCLUSION

In empirical work using data structures with time and cross-sectional dimensions, it is often difficult to estimate correlations in the errors accurately. The solution of setting these correlations to zero may not, however, be the wisest course. This paper suggests a technique whereby estimation can exploit correlations in the error structure across cross-sectional units. By imposing a factor structure, which can have an appealing economic interpretation, we provide a full-rank estimator of the covariance matrix. Monte Carlo shows that the parameter estimator thus obtained has substantial gains in efficiency over the standard alternative procedures for the data-sets considered, as well as offering unbiased inference. Our suggestion is that the proposed factor-residual method of estimation provides an important generalisation of standard panel techniques, and that even in situations where the conventional estimator of the covariance matrix of the residuals has full rank, it can provide a more parsimonious and efficient estimation technique.

APPENDIX A. APPENDIX

Properties of the Likelihood Function of the Factor Model

Our aim in this Appendix is to give a description of the geometry of the likelihood surface of the factor model, paying particular attention to the case when S is rank deficient. When S is of full rank it is possible to show by completely elementary methods that the likelihood function is bounded (Theorem 2 (i)). We do not have a similar simple demonstration that the maximum is attained (perhaps on the boundary). This is true, however, and can be shown by use of the concentration of L introduced by Lawley. When S is not of full rank, it remains true that L is bounded and attains its maximum provided that not too many factors are fitted, but here the tricky part is to show L is bounded. We achieve this by studying the behaviour of Lawley's concentrated L on the *rim*, the set of Σ s maximising L along rays from the origin.

Definitions and Simple Results

Let matrix $S = X'X$ where X is $T \times n$. Call S *admissible* if no column of X consists entirely of 0s. Define the function

$$(7) \quad L = L(\Lambda, \Psi, S) = -(\log \det \Sigma + \text{tr} \Sigma^{-1} S)$$

where Λ is any $n \times m$ matrix, Ψ is $n \times n$ diagonal with entries $\psi_i > 0$ and $\Sigma(\Lambda, \Psi) = \Lambda\Lambda' + \Psi$. Note Σ is the sum of a positive definite matrix and a non-negative definite matrix and is thus itself positive definite and hence invertible. As defined, L is the likelihood function (divided by $T/2$) of the m -factor model generated by n normal variates.

Some more terminology. Denote by \mathcal{L} the set of $n \times m$ matrices, by \mathcal{P} the set of $n \times n$ positive diagonal matrices, and by \mathcal{Q} the range of the function $\Sigma(.,.)$ so that $\Sigma : \mathcal{L} \times \mathcal{P} \rightarrow \mathcal{Q}$. This function is not one-to-one which means that Λ, Ψ are not identified by knowledge of $\Sigma(\Lambda, \Psi)$. In particular $\Sigma(\Lambda U, \Psi) = \Sigma(\Lambda, \Psi)$ for any unitary U . If required, this problem can be solved by restrictions on \mathcal{L} e.g. that it consist of Λ s which are non-negative and decreasing along the diagonal, and zero above it. So restricted, the function $\Sigma(.,.)$ is bi-continuous. This is not an important issue for our purposes as we seek the Σ s rather than the Λ s and Ψ s which generate them.

If D is a fixed $n \times n$ diagonal matrix then, up to a constant depending on D ,

$$L(\Lambda, \Psi, S) = L(D'\Lambda, D'\Psi D, D'SD)$$

Thus, for the study of boundedness and the existence of maxima, no generality is lost by re-scaling the rows and columns of S . When S is admissible, the elements along the main diagonal are non-zero and S may be re-scaled to a correlation matrix. This is often convenient.

A result we shall not need, but gives some insight into the character of the likelihood surface is:

THEOREM 1. The function L has no local minima.

PROOF. If Λ_0, Ψ_0 were a local minimum then L restricted to $\sqrt{\delta} \Lambda_0, \delta \Psi_0, \delta > 0$ has a local minimum for $\delta = 1$. Examination of $L(\delta) = L(\sqrt{\delta} \Lambda_0, \delta \Psi_0)$ shows this function has no local minimum. \square

It is well-known that L can have multiple maxima (e.g. Jöreskog, 1967). Thus likelihood surfaces can have ravines but, according to Theorem 1, no lakes.

We now set out the boundedness properties of the likelihood function.

PROPOSITION 1. If S is inadmissible, L is unbounded.

PROOF. If S is inadmissible then the j^{th} (say) row and column consist of 0s. Choose $\Lambda = 0$, so that $L = -(\sum_{i=1}^n \log \psi_i + \sum_{i=1}^n s_{ii}/\psi_i)$ and let $\psi_j \rightarrow 0$. There is a term in L , $-\log \psi_j$, which is unmatched by a term $-1/\psi_j$ so that $L \rightarrow \infty$. \square

PROPOSITION 2. If $\text{rank}(S) = n$, then L is bounded.

PROOF. Choose unitary U to diagonalise $\Sigma(\Lambda, \Psi)$ so that

$$L = - \sum_{i=1}^n (\log \sigma_i + x_i' S x_i / \sigma_i)$$

where $\sigma_1 \geq \sigma_2 \dots \geq \sigma_n > 0$ are the eigenvalues of Σ with corresponding eigenvectors x_i . Since $x' S x$ takes values in $[s_n, s_1]$ (where s_i denotes the eigenvalues of S and $s_1 \geq s_2 \dots \geq s_n > 0$) for vectors x with $\|x\| = 1$,

$$L \leq - \sum_{i=1}^n (\log \sigma_i + s_n / \sigma_i) \leq -n (\log s_n + 1)$$

because $-(\log \sigma_i + s_n / \sigma_i)$ is maximised at $\sigma_i = s_n$. \square

PROPOSITION 3. If $\text{rank}(S) = r < n$ and $m = r$ then L is unbounded.

PROOF. Consider, for $\mu > 0$,

$$\Sigma_\mu = \mu I + x_1 x_1' + \dots + x_r x_r'$$

where $r = \text{rank } S$ and $x_i, i = 1, \dots, r$ are the normalised eigenvectors of S of non-zero eigenvalues s_i . It follows that $\Sigma_\mu \in \mathcal{Q}$. Now the eigenvectors of Σ are precisely those of S with eigenvalues $\mu + 1$ for $i = 1, \dots, r$, and μ otherwise. It follows that we may simultaneously diagonalise Σ_μ and S in (7) to obtain

$$\begin{aligned} L &= - \left(r \log (1 + \mu) + (n - r) \log \mu + \sum_{i=1}^r s_i / (1 + \mu) \right) \\ &\rightarrow \infty \text{ as } \mu \rightarrow 0. \square \end{aligned}$$

If we were to fit m factors in Proposition 3 where $m \leq r < n$ then the likelihood function takes the form

$$L = - \left(m \log (1 + \mu) + (n - m) \log \mu + \sum_{i=1}^m \frac{s_i}{1 + \mu} + \sum_{i=m+1}^r \frac{s_i}{\mu} \right)$$

and it is easy to see that this function is bounded as a function of μ as long as $m < r$. This might lead one to conjecture that L is bounded provided one fits fewer than r factors. However this is incorrect: certain null space structures of the matrix S further reduce the number of possible factors. To see this suppose that for all $z \in \mathcal{N}(S)$ (the null space of S), $z_j = 0$ for some index j , $1 \leq j \leq n$, i.e. the null space has a row of zeros at the j^{th} position.¹⁶ This implies that $e_j \perp \mathcal{N}(S)$ where e_j is the j^{th} element of the canonical basis of \mathbb{R}^n . It follows that e_j is a linear combination of the eigenvectors of S of non-zero eigenvalue whence

$$\sum_{i=1}^r x_i x_i' + \mu I = \sum_{i=1}^{r-1} y_i y_i' + \Psi$$

where the y_i are linear combinations of the x_i , and Ψ is diagonal. This enables the likelihood for the r -factor model to be written as an $r - 1$ factor model. Since for $\text{rank}(S) = r$, the r -factor model gives an unbounded likelihood, each row of zeros in the null space reduces by one the number of factors that can be fitted. Thus the structure of the null space plays a key part in determining the number of factors that can be fitted. It turns out that the appropriate condition is a generalisation of the number of rows of zero in the null space which we now develop. For convenience we state the main result.

THEOREM 2. If S is inadmissible then L is unbounded. If S is admissible:

- (i) If $\text{rank}(S) = n$ then L is bounded
 - (ii) If $\text{rank}(S) = r < n$ and $m \geq r - d(S)$ then L is unbounded
 - (iii) If $\text{rank}(S) = r < n$ and $m < r - d(S)$ then L is bounded
- where $d(S)$ is the *defect* of matrix S , defined below.

Propositions 1-3 establish the first part. We now proceed to the proof of parts (ii) and (iii). This builds on the concentration of the likelihood introduced by Lawley.

Lawley's Machinery

The following results are essentially due to Lawley (1940, 1942, 1943); Lawley and Maxwell (1963) give a convenient condensed version; Jöreskog (1967) also provides a useful account:

Assume $\Sigma = \Sigma(\Lambda, \Psi)$. Then, by routine calculation,

$$(8) \quad \partial L / \partial \Lambda = -2\Sigma^{-1}(\Sigma - S)\Sigma^{-1}\Lambda$$

$$(9) \quad \partial L / \partial \Psi = -\text{diag}(\Sigma^{-1}(\Sigma - S)\Sigma^{-1})$$

$$(10) \quad \Sigma^{-1} = \Psi^{-1} - \Psi^{-1}\Lambda(I_m + \Lambda'\Psi^{-1}\Lambda)^{-1}\Lambda'\Psi^{-1}$$

If in addition, $\partial L / \partial \Lambda = 0$, then a little algebra shows

$$(11) \quad S\Psi^{-1}\Lambda(I_m + \Lambda'\Psi^{-1}\Lambda)^{-1} = \Lambda$$

$$(12) \quad \Sigma^{-1}(\Sigma - S)\Sigma^{-1} = \Psi^{-1}(\Sigma - S)\Psi^{-1}$$

$$(13) \quad \partial L / \partial \Psi = -\Psi^{-1}\text{diag}(\Sigma - S)\Psi^{-1}$$

¹⁶An abuse of terminology that we will often find, as here, too convenient to resist is to identify a linear space with its basis vectors written as a matrix.

We make the normalisations $S^* = \Psi^{-\frac{1}{2}} S \Psi^{-\frac{1}{2}}$, $\Lambda^* = \Psi^{-\frac{1}{2}} \Lambda$. Then (11) is transformed to

$$(14) \quad S^* \Lambda^* (I_m + \Lambda^* \Lambda^{*\prime})^{-1} = \Lambda^*$$

Equation (14) is essentially a collection of eigenvector equations from which optimal Λ^* can be obtained for each value of Ψ . Lawley-Jöreskog base an ML procedure on (14), searching over Ψ . This is for the full-rank case wherein the existence of an ML solution is guaranteed according to Proposition 3. In the rank-deficient case, with no such assurance in general, we need to proceed with a little more care. Nevertheless, we shall see below that, for given Ψ , optimal Λ are indeed determined by (14)

Concentrating out Λ

Substitution from (10) into (7) and a little manipulation yields

$$(15) \quad L = - \left[\log \det \Psi + \log \det (I_m + \Lambda^{*\prime} \Lambda^*) + \text{tr} S^* - \text{tr} \Lambda^{*\prime} S^* \Lambda^* (I_m + \Lambda^{*\prime} \Lambda^*)^{-1} \right]$$

Since $L(\Psi, \Lambda U) = L(\Psi, \Lambda)$ for any conformable orthogonal matrix U it follows that we may replace Λ^* by $\Lambda^* U$ in (15), choosing U so that the columns of $\Lambda^* U$ are orthogonal vectors or 0s. Assume this has been done and extract the terms in Λ^* from (15):

$$L_0 = - \sum_{i=1}^m [\log(1 + x_i' x_i) - x_i' S^* x_i / (1 + x_i' x_i)]$$

where the x_i are the columns of Λ^* . We wish to maximise L_0 over all systems of vectors $x_i, i = 1, \dots, m$ where the x_i are zero or pairwise orthogonal. Hold each $x_i' x_i$ fixed and regard L_0 as a function of the normalised x_i . Since the supremum of $\sum_{i=1}^m x_i' S^* x_i$ is attained at eigenvectors of S^* or zero vectors (Rao, page 63) it follows that L_0 is a sum of terms of the form: $s_i^* x_i' x_i / (1 + x_i' x_i) - \log(1 + x_i' x_i)$. These terms contribute non-negatively if and only if $s_i^* \geq 1$. If $s_i^* < 1$ then L_0 is maximised by choosing $x_i = 0$. The optimal system of $x_i, i = 1, \dots, m$ then consists of the set of eigenvectors of S^* with modulus determined by

$$(16) \quad 1 + x_i' x_i = \max(1, s_i^*)$$

Let $\pi(S^*)$ be the number of eigenvalues of S^* greater than 1 and define $m_0 = \min(\pi(S^*), m)$. The columns of the optimal $\Lambda^* = \Lambda^*(\Psi)$ in (15) thus consist of the first m_0 eigenvectors of S^* , with modulus determined by (16), 0s elsewhere. Define $\Lambda(\Psi) = \Psi^{\frac{1}{2}} \Lambda^*(\Psi)$. The concentrated likelihood function is now $L^c(\Psi) = L(\Psi, \Lambda(\Psi))$. This function is well-defined for all S , be it full-rank or rank-deficient, admissible or inadmissible.

Substituting in (15), one finds, up to a constant,

$$(17) \quad L^c(\Psi) = - \left[\log \det \Psi + \sum_{i=1}^{m_0} (\log s_i^* - s_i^*) + \text{tr} S^* \right]$$

This is the function we need to bound.

Define $\mathcal{F} = \{\Psi \in \mathcal{P}; \psi_i \leq s_{ii}\}$. Then since $\partial L/\partial \Lambda = 0$ at $(\Psi, \Lambda(\Psi))$ and $\partial L^c/\partial \Psi = \partial L/\partial \Psi$ (the envelope theorem) it follows from (13) that $\partial L^c/\partial \psi_i < 0$ for $\psi_i > s_{ii}$ with the implication that, for each $\Psi \notin \mathcal{F}$, there is an element of \mathcal{F} for which the concentrated L takes a higher value. We define the *rim* \mathcal{R} as the set of Ψ which maximise the concentrated likelihood function along rays from the origin. Although outside \mathcal{F} all directional derivatives are strictly negative, the rim does not necessarily lie within the bounded set \mathcal{F} ; however, given a point on the rim not within \mathcal{F} , one can find a dominating point within \mathcal{F} by moving along a directional derivative towards \mathcal{F} , subsequently passing back to the rim along a ray from the origin. It is easy to see that iterating this procedure leads to a dominating value in $\mathcal{R} \cap \mathcal{F}$. It follows that the maximum is attained on $\mathcal{R} \cap \mathcal{F}$ or on the boundary of \mathcal{F} . Thus we have:

THEOREM 3. If L is bounded, it attains its maximum in \mathcal{F} or $\overline{\mathcal{F}} \setminus \mathcal{F}$ in the sense that

$$\sup L = \lim_{p \rightarrow \infty} L^c(\Psi^p)$$

where $\lim_{p \rightarrow \infty} \Psi^p = \Psi^0 \in \overline{\mathcal{F}} \setminus \mathcal{F}$ and $\Psi^p \in \mathcal{R} \cap \mathcal{F}$ for all p .

Such boundary solutions where some $\psi_i = 0$ can occur and are called Heywood solutions. At a Heywood solution one need not have $\partial L/\partial \Psi = 0$. Since the likelihood function is evidently continuous where defined, it can be unbounded only at such boundary solutions. The study of the boundedness of the likelihood function thus consists in large part of the study of Heywood solutions. Examination of (17) reveals we need an apparatus to analyse the simultaneous behaviour of the ψ_i and the s_j^* as some of the $\psi_i \rightarrow 0$. We proceed to this.

The box diagram

To complete the proof of Theorem 2 we consider the spectrum of S^* . Assume S is admissible. Consider a sequence $\Psi^p \in \mathcal{F}$ with $\lim_{p \rightarrow \infty} \Psi^p = \Psi^0$ where $\psi_i^0 = 0$ for some i . By passing to a subsequence and renaming indices if necessary, we can assume

$$\psi_1^p \leq \psi_2^p \leq \dots \leq \psi_n^p$$

for each p . We assume $\psi_j^p \rightarrow 0$ if and only if $j < N$. Define S^{p*} corresponding to Ψ^p and assume it has eigenvalues $s_1^{p*} \geq s_2^{p*} \geq \dots \geq s_n^{p*} \geq 0$ with corresponding unit eigenvectors x_i^p . If S is rank-deficient; $r = \text{rank} S < n$, then the last $n - r$ of the s_i^{p*} are zero. Define $\Psi_i^p = s_i^{p*} \Psi^p$. Then the eigenvalue equation for S^{p*} takes the form

$$(18) \quad S (\Psi_i^p)^{-\frac{1}{2}} x_i^p = (\Psi_i^p)^{\frac{1}{2}} x_i^p \quad 1 = 1, \dots, r$$

For each p , the system $x_i^p, i = 1, \dots, r$, constitutes an orthonormal set in \mathbb{R}^n and the compactness of the unit ball implies there exists a limit orthonormal system $x_i^0, i = 1, \dots, r$ where $x_i^0 = \lim_{p \rightarrow \infty} x_i^p$ (passing to a subsequence if necessary). Treating (18) as an equation of the form $Sx = y$, one deduces

$$(\Psi_i^p)^{-\frac{1}{2}} x_i^p = S^{-1} (\Psi_i^p)^{\frac{1}{2}} x_i^p + z$$

where $z \in \mathcal{N}(S)$, and S^- is any generalised inverse of S , which we may take to be positive definite. Since $\mathcal{N}(S^{p*}) = (\Psi_i^p)^{\frac{1}{2}} \mathcal{N}(S)$, it follows that

$$(19) \quad x_i^p = (\Psi_i^p)^{\frac{1}{2}} S^- (\Psi_i^p)^{\frac{1}{2}} x_i^p + z^*$$

where $z^* \in \mathcal{N}(S^{p*})$. Thus, since $x_i^p \perp z^*$,

$$(20) \quad 1 = x_i^{p'} (\Psi_i^p)^{\frac{1}{2}} S^- (\Psi_i^p)^{\frac{1}{2}} x_i^p \quad i = 1, \dots, r$$

It follows from (20) that

$$(21) \quad 1/s_{\max}^- \leq \left\| (\Psi_i^p)^{\frac{1}{2}} x_i^p \right\|^2 \leq 1/s_{\min}^-$$

where s_{\max}^- and s_{\min}^- denote the largest and smallest eigenvalues of S^- , respectively. Thus, passing if necessary to a subsequence of Ψ^p , we deduce that $\psi_j^p s_i^{p*}$ approaches a finite limit on the support of x_i^0 (the indices r for which $x_{i,r}^0 \neq 0$) and a non-zero limit for at least one r in the support. The limiting behaviour of $\psi_j^p s_i^{p*}$ is indicated in the diagram:

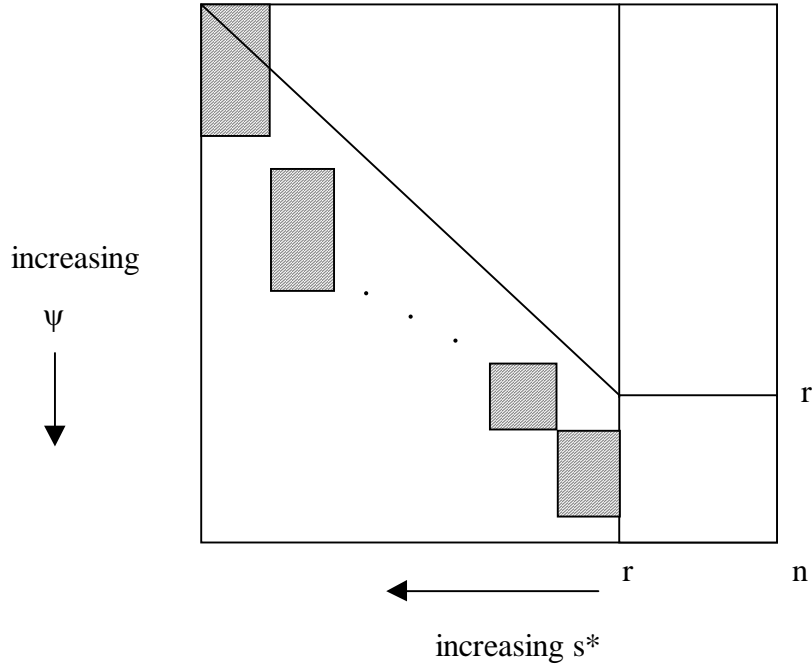


Figure A1 The Box Diagram

The shaded areas indicate (i, j) pairs for which $\psi_j s_i^*$ approaches a non-zero limit. The diagram incorporates the following properties:

1. The pairs (i, j) are arranged as non-overlapping boxes. Each $i \leq r$ corresponds to one box. If (i, j) and (i', j') belong to the same box then $\psi_j s_i^* / \psi_{j'} s_{i'}^*$, $\psi_j / \psi_{j'}$ and $s_i^* / s_{i'}^*$ all approach non-zero limits. To see this, observe that the boxes could be constructed as follows. For each i , the j for which $\psi_j s_i^*$ approaches

a non-zero limit form a vertical line-segment by the monotonicity of ψ_j with respect to j . Next note that if two such line-segments were to overlap (share a common j), they would lie in a box with the stated properties.

2. Above the boxes, $\psi_j s_i^* \rightarrow 0$. This follows by construction.
3. Beneath the boxes, $x_{ik}^0 = 0$. This follows by construction since $\psi_j s_i^*$ approaches a finite limit on the support of x_i^0 .
4. The boxes start in the top left hand corner. If not we would have $\psi_1 s_1^* \rightarrow 0$ and hence $\psi_1 s_i^* \rightarrow 0$ for all i . But $\text{tr} S^* = \sum_{i=1}^r s_i^* = \sum_{i=1}^n s_{ii}/\psi_i$ so then $0 = \lim \psi_1 \sum_{i=1}^r s_i^* \geq s_{11} > 0$.
5. The boxes are tall: their height is not exceeded by their breadth. To see this, first define

$$G(x_i^0) = \lim_{p \rightarrow \infty} \Psi^{p\frac{1}{2}} x_i^0$$

and extend to the span by linearity; define also

$$H(x) = (S^-)^{\frac{1}{2}} G(x).$$

If u, v on the horizontal axis belong to the same box then

$$\begin{aligned} H(x_u^0)' H(x_v^0) &= \lim_{p \rightarrow \infty} x_u^0' \Psi^{p\frac{1}{2}} S^- \Psi^{p\frac{1}{2}} x_v^0 \\ &\neq 0 \end{aligned}$$

since the limits of Ψ_u and Ψ_v differ by a non-zero multiplicative constant in the same box. It follows that H and hence G preserve dimension on the span of the x_i^0 corresponding to each box. But the support of each $G(x_i^0)$ lies within such a box so the result follows.

The Null Space

So far we have not defined limit eigenvectors in the null space $\mathcal{N}(S^{p*})$. Let z_{r+1}, \dots, z_n be a basis for $\mathcal{N}(S)$ in reduced column-echelon form and, for $j = 1, \dots, n$, define $V(j)$ as the largest index $r+k$ for which the leading 1 has row-number greater than or equal to j . In the event that the basis has rows of zeros at the bottom, we define $V(j) = r$ for these rows. The step-shaped path $(j, V(j))$ thus defines the frontier between the zero and non-zero regions of $\mathcal{N}(S)$. Let s be the minimum number of non-zero entries among all non-zero vectors of $\mathcal{N}(S)$ and define the *defect* of S , $d(S)$ by

$$s = r - d(S) + 1$$

Reduction to echelon form will produce $r + 1$ or fewer non-zero entries in z_n and $d(S) > 0$ indicates extra structure in $\mathcal{N}(S)$. This can be so only if an algebraic identity holds among the variates used to generate the covariance matrix S , typically a zero-probability event. Note that each row of zeros in the null space contributes one to the defect. It is possible, perhaps having first permuted the indices i , to write a basis for $\mathcal{N}(S)$ in column-echelon form so that $V(r - d(S) + 1) = n$. In this case the right-most column in the null space has non-zero entries for $j \leq r - d(S) + 1$, zeros elsewhere. In general (without necessarily permuting the indices), column-echelon form delivers $V(r - d + 1) = n$ for some $d \geq d(S)$ where $r - d + 1$ is the number of terms lying above (and including) the leading 1 in z_n .¹⁷

Column-echelon form ensures the following properties of the function V :

1. $V(j) = n$ for $j \leq r - d + 1$.

¹⁷As defined, d is invariant to the particular column-echelon form.

2. $0 \geq \Delta V(j) \geq -1$, where $\Delta V(j) = V(j) - V(j-1)$.
3. $\Delta V(r-d+1) = -1$.

Now define $F(j)$ by

$$V(j) - F(j) = n - j \text{ for } j = 1, \dots, n$$

Then F inherits the properties:

1. $F(j) = j$ for $j \leq r-d+1$.
2. $1 \geq \Delta F(j) \geq 0$.
3. $\Delta F(r-d+1) = 0$.
4. $i - F(j) \leq 1$ for $j > r-d+1$,

where property 4 follows from 2 and 3 immediately above.

Limit eigenvectors in $\mathcal{N}(S^{p*})$ are constructed as follows. In the event that $N \leq r-d+1$ then $x_{r+i}^0 = \Psi^{0\frac{1}{2}} z_{r+i}$, $i = 1, \dots, n-r$, are linearly independent, each a limit of the sequence $\Psi^{p\frac{1}{2}} z_{r+i}$ and orthogonal to x_i^0 , $i = 1, \dots, r$. When $N > r-d+1$, the last $n - V(N)$ vectors as defined above vanish, but these can be replaced by $\lim_{p \rightarrow \infty} (\Psi^p / \psi_k)^{\frac{1}{2}} z_{r+i}$, where k is the row number of the leading 1 in z_{r+i} . These are all limits of sequences in $\mathcal{N}(S^{p*})$ and have some non-zero elements for indices j less than N . In both cases we are led to limit null-space vectors x_{r+1}^0, \dots, x_n^0 for which: (a) If $N < r-d+1$ then the limit vectors have zeros for indices less than N . (b) If $N \geq r-d+1$ then the limit vectors can be divided into two groups, one with zero elements at indices above N , one with zero elements at indices below N , the dividing vertex being $(N, V(N))$. This structure is summarised in the Figure below.

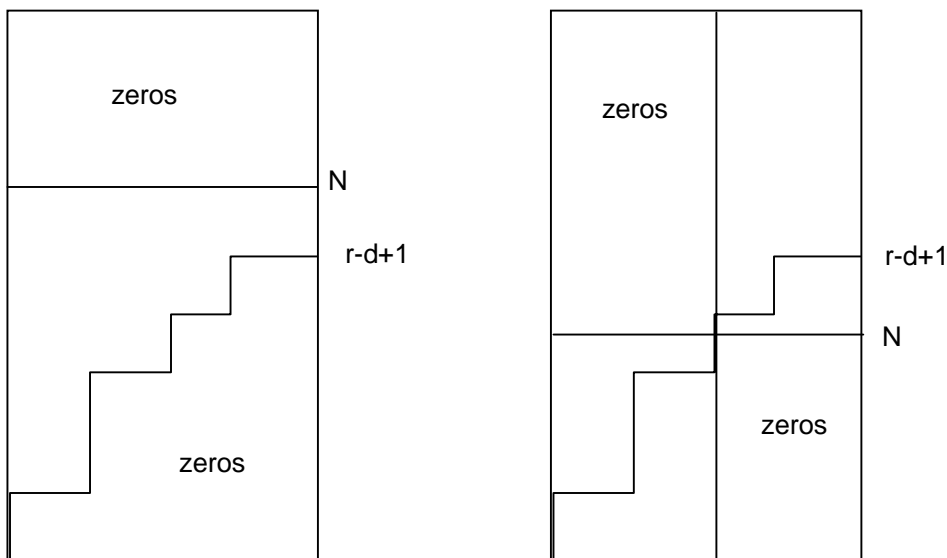


Figure A2 Location of zeros in the limit null space

LEMMA 1. In the box-diagram, the right-most box lies beneath $j = N$ if and only if $F(N) < r + 1$. In this case, the North-West vertex of the box lies on the path $(N, F(N))$.

PROOF. We prove the second part first. Assume therefore that the right-most box lies beneath $j = N$. Taking limits in (19), we deduce that the subvectors of each x_i^0 , $i = 1, \dots, r$ lying above $j = N$ belong to the limit null-space, analogously truncated. Exploiting the orthogonality structure implicit in Figure 2, and the mutual orthogonality between x_1^0, \dots, x_r^0 and x_{r+1}^0, \dots, x_n^0 , one can deduce that these subvectors are identically zero. It follows that the limit eigenvectors in the right-most box, together with the limit null space vectors up to index $V(N)$, all have zeros above $j = N$, whereas all other eigenvectors have zeros below $j = N$. Figure 3 illustrates the argument (note we have drawn the V and F functions as straight lines for simplicity). Considerations of dimensionality now show that the submatrix of the right-most box lying beneath $j = N$ shown shaded in Figure 3, together with the correspondingly truncated vectors in the adjacent region of the limit null space, forms a square. The result now follows from the definition of the function F .

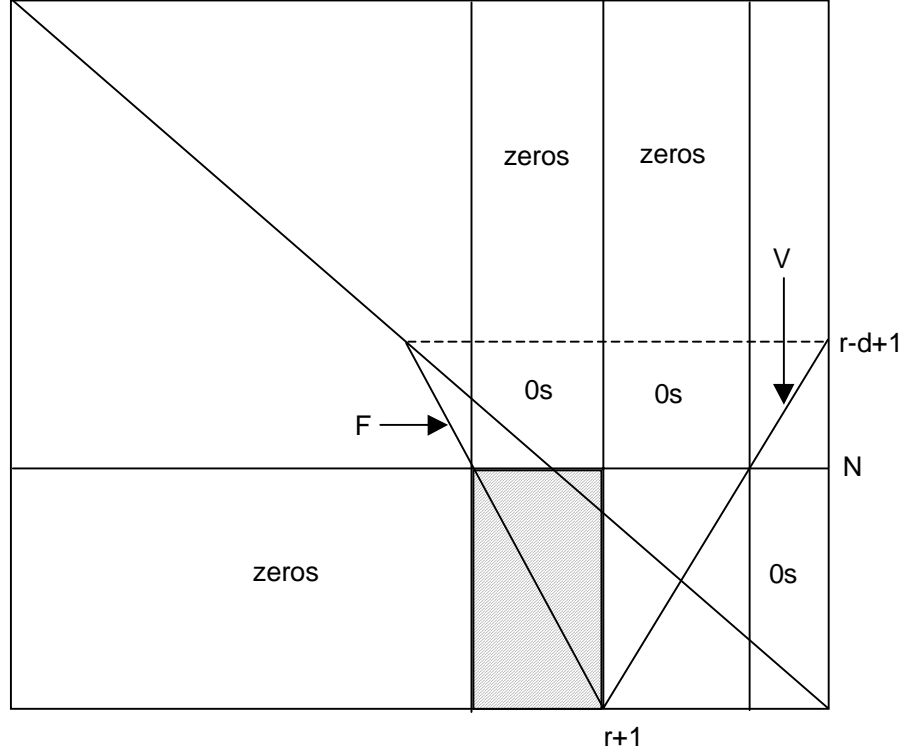


Figure A3 Interaction between the null space and the right-most box

If there is no box beneath $j = N$ then the above argument shows that the first block of truncated vectors in the null space themselves form a square, and are linearly independent by construction. It follows that

$$V(j) = j + r \text{ for } j = 1, \dots, N - r$$

whence $F(j) = 1 + r$. Clearly this argument is reversible, so the Lemma is proved. \square

Proof of Theorem 2

We first prove boundedness for $m < r - d(S)$. Assume, by way of contradiction, L^c is unbounded for some sequence $\Psi^p \rightarrow \Psi^0$. In choosing the elements ψ_j^p to be monotonic in j , we have implicitly re-ordered the indices j . We have however $m < r - d(S) \leq r - d$ where $V(r - d + 1) = n$. Clearly we may take all Ψ^p to lie in $\mathcal{R} \cap \mathcal{F}$ so that $\Psi^0 \in \overline{\mathcal{F}} \setminus \mathcal{F}$. We first prove the following.

LEMMA 2. On the rim,

$$(22) \quad tr S^* - \sum_{i=1}^{m_0} s_i^* = n - m_0$$

and, up to an additive constant,

$$(23) \quad L^c = - \left(\sum_{i=1}^n \log \psi_i + \sum_{i=1}^{m_0} \log s_i^* \right)$$

Moreover, $m = m_0$, i.e. $s_i^* \geq 1$ for $i \leq m$.

PROOF. For $\delta > 0$, one has from (17)

$$L^c(\delta\Psi) = - \left(\log \det \Psi + \sum_{i=1}^{m_0} \log s_i^* + (n - m_0) \log \delta + (tr S^* - \sum_{i=1}^{m_0} s_i^*)/\delta \right)$$

Now $tr S^* = \sum_{i=1}^n s_i^*$ so $tr S^* > \sum_{i=1}^{m_0} s_i^*$ in virtue of the assumption that $m < \text{rank } S = \text{rank } S^*$. The integer m_0 can depend on δ but, irregardless, $L^c(\delta\Psi)$ can be made arbitrarily negative by choice of δ sufficiently small or large. Continuity of $L^c(\delta\Psi)$ then guarantees that a maximum, at δ_m say, is attained. Fixing m_0 to be the value at this maximum, one deduces that $\delta_m = (tr S^* - \sum_{i=1}^{m_0} s_i^*)/(n - m_0)$.

Since $\delta_m = 1$ if Ψ already belongs to the rim, the first result follows and the second is obtained by direct substitution in 17. Finally, m_0 can be smaller than m only if $s_{m_0+1} < 1$; however 22 implies that the average of the numbers $s_i^*, i = m_0 + 1, \dots, n$, some of which are zero, is unity, inconsistent with this. \square

The implication of (22) is that, for Ψ^p on the rim, $0 \leq s_i^* \leq 1$ for $i = m + 1, \dots, n$. We may thus choose a subsequence of Ψ^p for which each s_i^* converges, which implies that $\psi_j s_i^*$ converges to a finite limit for $j \geq N$, $r \geq i > m$. (22) further implies that at least one of these limits is non-zero. It follows that the right-most box lies beneath $j = N$ in the box diagram and hence that the North-West vertex of this box is given by $(N, F(N))$ where $m + 1 \geq F(N)$. If $N > r - d + 1$ then, since $m < r - d$, we have

$$m + 1 < r - d + 1 \leq F(N)$$

by properties (i) and (ii) of F above. This is contrary to $m + 1 \geq F(N)$ and we conclude that $N \leq r - d + 1$.¹⁸ It follows that $F(N) = N$ by property (i) of F above i.e. the NW vertex of the right-most box lies on the diagonal. But since boxes are "tall", all boxes must then lie along the diagonal with the implication that $\psi_i s_i^*$ converges to a non-zero limit for all $i \leq r$. If (23) is recast in the form

$$L^c = - \left(\sum_{i=1}^m \log \psi_i s_i^* + \sum_{i=m+1}^n \log \psi_i \right)$$

then boundedness now follows since $N = F(N) \leq m + 1$. This completes the proof of Theorem 2(iii).

Now assume that $m \geq r - d(S)$, $r < n$. We assume $d > 0$, for otherwise unboundedness would follow from Proposition 3. Thus $m \leq r - 2$. We re-order the indices j so that a reduced column-echelon form for $\mathcal{N}(S)$ has leading 1 in place $r - d(S) + 1$ for z_n . We choose a convergent sequence Ψ^p with $N = m + 2$ so that $n \geq N > r - d(S) + 1$. No harm is done by assuming m_0 is fixed as p changes. If

¹⁸This argument has shown, in fact, that $N > r - d + 1$ implies the sequence Ψ^p is not on the rim.

$m_0 < m$, we must have $s_{m_0+1}^{p*} \leq 1$, so it follows that $m_0 + 1 \geq F(N)$. Write (17) in the form

$$L^c = - \left(\sum_{i=1}^{m_0} \log \psi_i s_i^* + \sum_{i=m_0+1}^n s_i^* + \sum_{i=m_0+1}^n \log \psi_i \right)$$

Considering these three terms in turn, note that the first approaches a finite limit or $+\infty$ since the diagonal in the box diagram intersects the boxes or lies above them; the second is bounded because since $m_0 + 1 \geq F(N)$; while the third contains a term in $-\log \psi_{m_0+1}$ which approaches $+\infty$ since $m_0 + 1 < N$. This proves Theorem 2(ii) so the proof of Theorem 2 is now complete.

Figure 4 gives the box-diagram for a convergent ML estimation of the factor model.

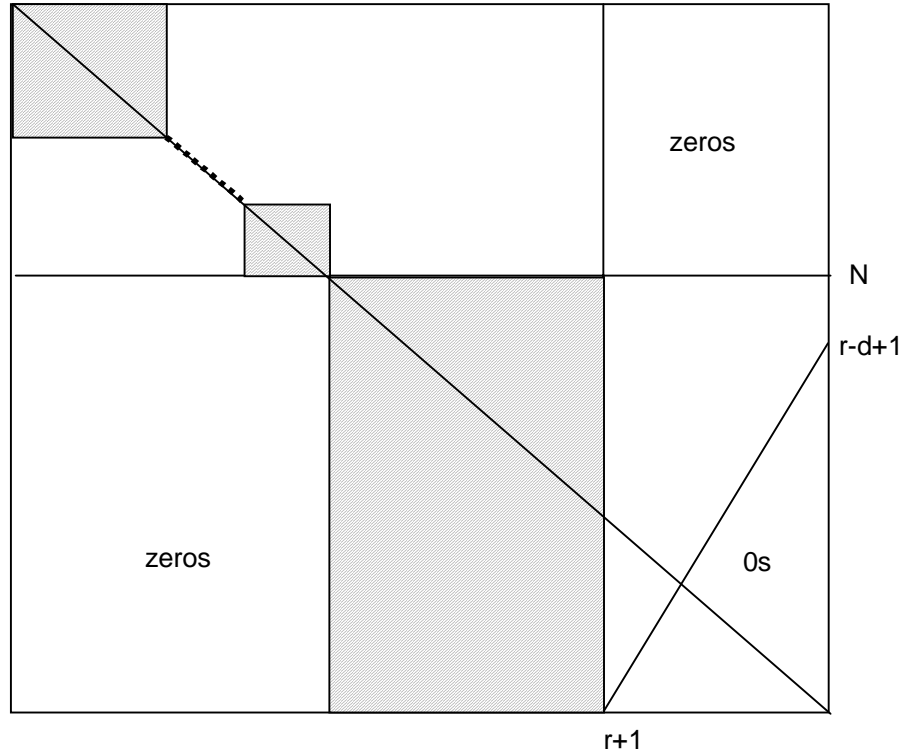


Figure A4 The box diagram for convergent maximum likelihood

In the shaded boxes, $\psi_j s_j^*$ approaches a non-zero limit, above them, zero; beneath the boxes the limit eigenvectors x_i^0 have zero support. The figure also gives the zero-structure for the limit null space. We have also established that $m \geq N - 1$. Define the multiplicity of a Heywood solution as the number of indices for which $\psi_j \rightarrow 0$, $N - 1$ in our terminology. Thus our proof of Theorem 2 has established:

COROLLARY. The multiplicity of a Heywood solution is less than or equal to the number of fitted factors.

THEOREM 4. Assume $m < r - d$. Then there exists an invertible Σ^0 in $\mathcal{R} \cap \mathcal{F}$ or its boundary to maximise L .

PROOF. By Theorems 2 and 3, if $m < r - d$, there exists a sequence Ψ^p in $\mathcal{R} \cap \mathcal{F}$ with limit Ψ^0 such that $\sup L = \lim_{p \rightarrow \infty} L^c(\Psi^p)$. It was shown in the proof of 2(iii) that $F(N) = N$. If Σ^0 corresponds to Ψ^0 then

$$\begin{aligned} \Sigma^0 &= \lim_{p \rightarrow \infty} (\Psi^p + \Lambda^p \Lambda^{p'}) \\ (24) \quad &= \lim_{p \rightarrow \infty} \Psi^{p\frac{1}{2}} (I + \Lambda^* \Lambda^{*'}) \Psi^{p\frac{1}{2}} \end{aligned}$$

where Λ^* consists of the first m eigenvectors of S^{p*} with modulus determined by (16). Let x_i^p , $i = 1, \dots, n$, be an orthonormal basis of eigenvectors of S^{*p} , ranked by eigenvalue, and define $q_i^p = s_i^{*p}$ for $i \leq m$, 1 for $i > m$. Then $I = \sum_{i=1}^n x_i^p x_i^{p'}$ so it follows that

$$\begin{aligned} \Sigma^0 &= \lim_{p \rightarrow \infty} \Psi^{p\frac{1}{2}} \left(\sum_{i=1}^n q_i^p x_i^p x_i^{p'} \right) \Psi^{p\frac{1}{2}} \\ &= \sum_{i=1}^{N-1} G(x_i^0) G(x_i^0)' + \sum_{i=N}^n q_i^0 \Psi^{0\frac{1}{2}} x_i^0 x_i^{0'} \Psi^{0\frac{1}{2}} \end{aligned}$$

where G is as defined in the proof of property (v) of the box-diagram, $\psi_j^0 > 0$ for $j \geq N$, and $q_i^p \geq 1$. The $G(x_i^0)$ are linearly independent in the same box and orthogonal between boxes, as well as orthogonal to each $\Psi^{0\frac{1}{2}} x_i^0$, $i \geq N$, in virtue of the structure of the box-diagram. Exploiting the fact that the $\Psi^{0\frac{1}{2}} x_i^0$ vanish above $j = N$ for $i \geq N$, one deduces that the $\Psi^{0\frac{1}{2}} x_i^0$ are linearly independent and the result follows. \square

Maximising the Likelihood

There is a voluminous literature on ML estimation of the factor model (see e.g. Lee and Jennrich, 1979, Rubin and Thayer, 1982). Analytic maximisation function is intractable, so numerical techniques are used. Our prime requirement for numerous Monte Carlo experiments was guaranteed convergence, at the expense, if necessary, of speed. Our reading of the literature on factor ML algorithms is that choice of algorithm can affect convergence speed by a factor of about 3, though not by factors of 10 or more. Given the speed of modern machines this would indicate that the particular algorithm chosen may be of secondary importance. After some experimentation, we have chosen a mixture of steepest ascent and Fletcher-Powell algorithms. Use of steepest ascent alone can be very time consuming, the geometry of the likelihood function being such that a large number of very small steps can be required (hem-stitching). Fletcher-Powell, on the other hand, can get stuck where, in the search direction, there is no appreciable likelihood value improvement, but the gradient is non-zero in a different direction i.e. the likelihood function is not well-approximated by an ellipsoid in certain regions. Mixing the two schemes ensures that the speed of Fletcher-Powell is used where it improves the likelihood function, but, if it runs into problems, we switch to the more robust steepest ascent to move the search along. We also implement a scheme whereby the step length increases following a successful step (i.e. one that improves the likelihood). It seems likely that the appropriate mix of steepest ascent and Fletcher-Powell and the step expansion, is problem dependent; we have not attempted to optimise in this direction. We have found that this mixture provides a very robust maximisation

procedure. For the regressions in this paper we used the following scheme. The starting values are derived by assuming that Ψ_0 is scalar with magnitude the average of the smallest $n - m + 1$ eigenvalues of S , and as Λ_0 the implied square root of $S - \Psi_0$. The search method uses steepest ascent for the first 100 steps, then switches to Fletcher-Powell. Following a successful step, the step-length is expanded by 20%. Convergence is checked by the change in the likelihood function being below some tolerance. At this point the gradient of the likelihood function is also checked. If the gradient exceeds a critical value we switch back to steepest ascent for a further 40 steps and then revert to Fletcher-Powell. This process continues until both the change in the likelihood and the slope of the likelihood fall beneath required thresholds. Convergence characteristics with this algorithm seem good; in several thousand Monte Carlo we have never failed to find a maximum.

REFERENCES

- [1] Barro, R. (1991) "Economic Growth in a Cross Section of Countries", *Journal of Political Economy*, 100(2) pp 223-51.
- [2] Caselli, F. Esquivel, G and F. Lefort, (1996) "Reopening the Convergence Debate: A New Look at Cross Country Growth Empirics" *Journal of Economic Growth*, September, pp 363-389.
- [3] Chamberlain, G and M.Rothschild (1983) "Arbitrage and Mean Variance Analysis on Large Asset Markets", *Econometrica*, 51, pp1281-1304.
- [4] Conley, T.G. (1999) "GMM estimation with cross sectional dependence", *Journal of Econometrics*, 92, pp1-45.
- [5] Connor, G and R.Korajczyk (1986) "Performance Measurement with the Arbitrage Pricing Theory: A New Framework for Analysis" *Journal of Financial Economics*, 15, pp373-394
- [6] Connor, G and R.Korajczyk (1988) "Risk and Return in an Equilibrium APT: Application of a New Test Methodology" *Journal of Financial Economics*, 21, pp255-289.
- [7] Di Liberto, A. and J.Symons. (1998) "Econometric Issues in Convergence Regressions" (mimeo).
- [8] Driscoll, J.C. and A.C. Kraay (1998) "Consistent Covariance Matrix Estimation with Spatially Dependent Panel Data", *The Review of Economics and Statistics*, 80, pp 549-560.
- [9] Evans, P. and G.Karras (1996) "Convergence Revisited", *Journal of Monetary Economics*, 37, pp 249-265.
- [10] Fisher, S. (1993) "The role of macroeconomic factors in growth", *Journal of Monetary Economics* pp 485-512.
- [11] Frees, E.W. (1995) "Assessing cross-sectional correlation in panel data", *Journal of Econometrics*, 69, pp393-414
- [12] Greene, W (1990) *Econometric Analysis*, Macmillan.
- [13] Islam, N (1995) "Growth Empirics: A Panel Data Approach" *Quarterly Journal of Economics*, November.
- [14] Lee, S.Y. and R.Jennrich (1979) "A Study of Algorithms for Covariance Structure. Analysis with Specific Comparisons Using Factor Analysis", *Psychometrika*, 44(1), pp 99-113.
- [15] Jöreskog, K.G. (1967) "Some Contributions to Maximum Likelihood Factor Analysis", *Psychometrika*, 32(4), pp443-482.
- [16] Keane, M. and D.Runkle (1992) "On the Estimation of Panel Data with Serial Correlation when Instruments are not Strictly Exogenous", *Journal of Business and Economic Statistics*, 10(1) pp 1-9.
- [17] Lawley, D.N. (1940) "The Estimation of Factor Loadings by the Method of Maximum Likelihood", *Proceedings of the Royal Society of Edinburgh Section A*, 60, pp 64-82
- [18] Lawley, D N. (1942) "Further Investigations in Factor Estimation" *Proceedings of the Royal Society of Edinburgh Section A*, 61, pp176-189
- [19] Lawley, D.N. (1943) "The application of the maximum likelihood method to factor analysis" *British Journal of Psychology*, 33, pp172-175
- [20] Lawley, D.N. and A.E.Maxwell, (1963) *Factor Analysis as a Statistical Method*, London, Butterworths
- [21] Liang, K-Y. and S.L.Zeger (1986) "Longitudinal Data Analysis using Generalised Linear Models" *Biometrika* 73,1, pp13-22.
- [22] Lee, K., Pesaran, H. and R.Smith (1995) "Growth and Convergence; A Multicountry Empirical Analysis of the Solow Growth Model" (mimeo)
- [23] Mankiw, N. Romer, D. and D.Weil (1992) "A Contribution to the Empirics of Economics Growth", *Quarterly Journal of Economics*, 107(2), pp 407-37.
- [24] Maxwell,A.E. (1981) "Factor Analysis" in *Encyclopedia of Statistical Sciences* Vol 3 ed. S.Kotz and N.L.Johnson, John Wiley & Sons.
- [25] Neyman, J and E.Scott, (1948) "Consistent Estimates Based on Partially Consistent Observations" *Econometrica*, 16(1) pp1-32.
- [26] Nickell, S.J. (1981) "Biases in Dynamic Models with Fixed Effects", *Econometrica*, 47(5), pp 1249-66.
- [27] Rao, C. R. (1973) *Linear Statistical Inference and its Applications*, Wiley, New York.
- [28] Rubin D and D.Thayer (1982) "EM Algorithms for ML Factor Analysis", *Psychometrika*, 47(1), pp 69-76.
- [29] Sargan, J. (1988) *Lectures on Advanced Econometric Theory*, Blackwell.

- [30] Theil, H. (1971) *Principles of Econometrics*, Wiley, New York
- [31] Zellner, A. (1962) "An Efficient Method of Estimating Seemingly Unrelated Regressions and Tests for Aggregation Bias" *Journal of the American Statistical Association*, vol 57, pp 348-368.

UNIVERSITY OF CAMBRIDGE

E-mail address: `donald.robertson@econ.cam.ac.uk`

UNIVERSITY COLLEGE LONDON

E-mail address: `james.symons@ucl.ac.uk`