# Empirical Limits for Time Series Econometric Models

Werner Ploberger        and        Peter C. B. Phillips
*University of Rochester*        *Cowles Foundation, Yale University*

4 January 1998

## Abstract

This paper seeks to characterize empirically achievable limits for time series econometric modeling. The approach involves the concept of minimal information loss in time series regression and the paper shows how to derive bounds that delimit the proximity of empirical measures to the true probability measure in models that are of econometric interest. The approach utilizes generally valid asymptotic expressions for Bayesian data densities and works from joint measures over the sample space and parameter space. A theorem due to Rissanen is modified so that it applies directly to probabilities about the relative likelihood (rather than averages), a new way of proving results of the Rissanen type is demonstrated, and the Rissanen theory is extended to nonstationary time series with unit roots, near unit roots and cointegration of unknown order. The corresponding bound for the minimal information loss in empirical work is shown not to be a constant, in general, but to be proportional to the logarithm of the determinant of the (possibility stochastic) Fisher-information matrix. In fact, the bound that determines proximity to the DGP is generally path dependent, and it depends specifically on the type as well as the number of regressors. Time trends are more costly than stochastic trends, which, in turn, are more costly than stationary regressors in achieving proximity to the true density. The conclusion is that, in a very real sense, the 'true' DGP is more elusive when there is nonstationarity in the data. Some implications of these results for prediction and for the achieving proximity to the optimal predictor are explored.

*Keywords:* Proximity Bounds, Data generating process, Empirical measures, Fisher information, Minimal information loss, Lebesgue measure, Optimal predictor, Path dependence, Trends, Unit roots,

*JEL Classification:* C22

---

# 1 Introduction

The objective of most statistical analysis, including studies in economic time series, is the construction of good empirical models for given data. The true model, or probability measure, for the data is unknown and, in most practical cases, it is reasonable to suppose that it is unknowable. It is usually hypothesized up to a parameter that needs to be estimated from the data. Often, the data is scarce relative to the number of parameters that need to be estimated, and this makes it intuitively evident that 'lower' dimensional parameter spaces may be preferable in practice to 'higher' dimensional ones, a maxim that governs much empirical work in statistics and econometrics.

The mathematical justification for this maxim of parsimony is important and is especially relevant in the context of models used in economic time series, where the series are often comparatively short. The present paper follows an approach pioneered by Rissanen (1986, 1987, 1996) in addressing this question and seeks to establish a theory of minimal information loss in time series regression that is suitable for use in modern econometrics settings. A survey of the field is given in Gerencser and Rissanen(1992) and the volume by Keuzenkamp, McAleer and Zellner (1999) contains papers that report on some recent developments. So far, Rissanen's ideas have had little impact in econometrics or on thinking about econometric methodology, although their importance was emphasized recently in Phillips (1996). Suppose data is available and all that is known is that the data-generating process (DGP) belongs to a $k$-dimensional parametric family and satisfies certain regularity conditions. The seminal theorem by Rissanen that we build on here shows that the minimum information distance (based on the relative likelihood) between any candidate probability measure and the true measure is, on average, bounded from below by the product of $k$ and the logarithm of the sample size for almost all parameters, i.e., all besides a Lebesgue null set. The bound provides a yardstick for how 'close' to the true probability measure we can get within a parametric family, assuming that the parameters all have to be estimated with the given data.

The present paper works with a broader class of assumptions than Rissanen, allowing for some nonstationary as well as stationary time series, so that the results apply to models with integrated and cointegrated variables as well as stationary time series. Some new techniques for proving results of this type are also developed here and these may be of some independent interest.

Results on minimal information loss turn out to have intimate connections with Bayesian modelling, and some of these connections are explored here. In particular, we show that Bayesian models (in the sense of Phillips and Ploberger, 1996 and Phillips, 1996) are 'nearly optimal' descriptions of the 'true' data generating process (DGP) given that the parameters are unknown.

The paper is organised as follows. Section 2 gives some modelling preliminaries, discussing both Bayesian and classical versions of empirical models. Our main result on the limits for empirical econometric models is contained in Section 3 and is derived under some high level assumptions, which are justified for some specific econometric

models in Sections 4 and 5. Section 6 explores some of the implications of our results for forecasting and achieving proximity to the optimal predictor . Section 7 concludes. Proofs and some complementary technical material are provided in an Appendix in Section 8.

## 2    Modelling Preliminaries

We start by considering a fairly typical empirical modelling situation with time series data. We have data $x^n = (x_t)_1^n$ that we associate with the realisation of a random process that takes values in a space $E$ with an associated event $\sigma-$ algebra $\mathfrak{F}$. The random elements need not be finite dimensional real vectors, and $E$ could be an arbitrary Polish space. So we can describe qualitative as well as quantitative data. The data are assumed to arrive consecutively, i.e., we get observation $x_n$ at 'time' $n$. We use $\mathfrak{F}_n$ to denote the information available at $n$ — i.e., $\mathfrak{F}_n \supset \sigma(x^n)$, the $\sigma$-algebra generated by $x^n$.

Our purpose is the evaluation of empirical models and, therefore, we need to clarify what we mean by this notion in a general context. Once the model concept is defined we have a natural basis for developing a criterion for relating different empirical models of the same process given the same observed data. In our framework, we think of a model as a sequence of conditional probability measures, $G_n$, from $\mathfrak{F}_n$ to $E$, i.e., a model is a representation of the process that allows us at each point $n$ and for every $x^n$ to calculate a prediction of the next observation in the random sequence. This is precisely what the conditional measure provides, viz., a mathematical description of a law that governs the forthcoming observation given the past that has been observed so far. Note that prediction is not taken here in the narrow sense of a linear prediction or projection on the past $x^n$, although it could turn out that this is one of its features. Instead, it is a complete probability distribution. It is easily seen that there is a one to one correspondence between models $(G_n)$ and measures $G$ on $E^{\mathbb{N}}$. In particular, due to the fact that $E$ is Polish, we can see that for every sequence $(G_n)$ it is possible to construct a compatible measure $G$, i.e., a measure whose conditional distributions are the $G_n$ and vice versa.

How do we find candidate empirical models of the data? There is some difference here between the stylized 'classical' and 'Bayesian' paradigms of data analysis. Our approach seeks to cover both paradigms. Let us assume that we are in a typical parametric context wherein the DGP is assumed to be known up to a certain parameter $\theta$ and let $P_\theta$ be the corresponding probability measure. The classical procedure is to use the information in $\mathfrak{F}_n$ to estimate $\theta$, say by the maximum likelihood estimator (MLE) $\hat{\theta}_n$, and then use $P_{\hat{\theta}_n}(\cdot|x^n)$ as the inferred empirical model for the process. One way of constructing an empirical measure from the classical framework is simply to 'plug' the estimator into the conditional probability measure in a recursive manner as we move through the data from some given point of initialization. In the terminology of Dawid (1984), the outcome of this recursion is a prequential density. On the other hand, in the Bayesian paradigm, a prior density $\pi(\theta)$ for $\theta$ is defined

and the Bayesian mixture

$$P = \int \pi\left(\theta\right) P_\theta d\theta \qquad (1)$$

gives the marginal distribution of the data $x^n$. We can then construct the conditional models (or measures) from $P$ and compute the associated data densities, viz.

$$pdf\left(x^n\right) = \frac{dP}{d\lambda} = \int \pi\left(\theta\right) \frac{dP_\theta}{d\lambda} d\theta, \quad pdf\left(x^n|x^{n_0}\right) = \frac{pdf\left(x^n\right)}{pdf\left(x^{n_0}\right)}, \qquad (2)$$

where $\lambda$ is a dominating measure (possibly Lebesgue measure) for $P_\theta$.

In the above setting, the class of potential models for the data is very wide. Indeed, as soon as we have a rule for obtaining numerical values of parameters or rules for averaging the parameters out, almost anything can be considered as a 'model' for the data. To prevent modelling concepts from degenerating into the trivial, we introduce a yardstick for measuring the 'goodness' of a model. Suppose the data are generated by some probability measure $P_\theta$ and we use a 'model' $G$ (i.e. a probability measure) as the supposed data generating mechanism. Denote by $P_\theta^{(n)}$ and $G^{(n)}$ the *restrictions* of these measures to $\mathfrak{F}_n$: i.e. we limit the information to that available at time $n$. Similarly, we denote by $P^{(n)}$ the restriction of the Bayesian measure $P$ to $\mathfrak{F}_n$. Then, our measure of 'goodness of fit' is just the sequence of random variables

$$\ell_n(G) = \log \frac{dG^{(n)}}{dP_\theta^{(n)}}.$$

These random variables allow us to compare different models (i.e., $G_1$ is 'better' than $G_2$ iff $\ell_n(G_1)$ 'is greater than' $\ell_n(G_2)$ — in whatever way we define an ordering between random variables), although the ordering is only a partial ordering because it is possible that some models are not compatible.

We think that this measure for the 'distance' of a given model from the 'true' probability measure is a sensible formalisation of the intuitive concept of one model being 'better' than another for the following reasons:

1. It is compatible with Kullback–Leibler (KL) type information 'metrics' since $-E_\theta \ell_n(G)$ is just the KL information distance of $G^{(n)}$ to $P_\theta^{(n)}$ (i.e. the measures modelling information up to time $n$). So if $G_1$ is better than $G_2$, then $G_1^{(n)}$ is, in KL-distance, nearer to $P_\theta^{(n)}$ than $G_2^{(n)}$.

2. If $\ell_n(G) = 0$, then $G^{(n)} = P_\theta^{(n)}$ , i.e., the measures describing the data are identical.

3. If, for $n \to \infty$, $\ell_n(G) = O_{P_\theta}(1)$ then $G^{(n)}$ and $P_\theta^{(n)}$ are *contiguous* in the sense of LeCam (1986). As a consequence, it is impossible to construct *consistent* tests of $P_\theta^{(n)}$ against $G^{(n)}$. So, in this case, it is impossible even asymptotically to tell for sure if the data were generated by $P_\theta$ or $G$.

4. Suppose we have given two models, say $G_1$ and $G_2$, and $\ell_n(G_1) - \ell_n(G_2) \to \infty$. If a researcher has to decide between these two models - i.e. choose the one which describes the data in a better way, then the Neyman–Pearson lemma suggests the use of the likelihood ratio (LR) test of $G_1$ against $G_2$. In this case, the researcher will choose the 'better' model (in our sense of the term) — since $\log \frac{dG_1}{dG_2} = \ell_n(G_1) - \ell_n(G_2) \to \infty$ asymptotically.

Having established the 'distance' between an empirical model and the 'true measure' and between one model and another, the question of finding the 'best' model arises naturally. In Phillips (1996) and Phillips and Ploberger (1996), the 'goodness' of Bayesian models was analysed in the context of a model with asymptotic locally quadratic likelihood, and a corresponding asymptotic approximation to the data density for such a model was computed. Those computations led to an empirical density for the data which was called a PIC density and this density was used for model selection purposes.

In this paper we generalize a result of Rissanen (1986, 1987), and we will show that the empirical PIC density is essentially optimal in terms of its rate of approximation to the true model. Given a parametrized family of probability measures and an empirical model for the data, we show that the Lebesgue-measure of the set of parameters corresponding to probability measures for which the model is 'better' than a certain bound converges to zero — i.e., the set for which we can beat this bound is relatively thin. Furthermore, the bound is shown to be achievable and is attained by the PIC density, clarifying the sense in which this model is optimal. A trivial example illustrating the sort of situation where the bound can be exceeded (i.e the thin set referred to above) is the empirical model consisting of a probability measure $G_n$ obtained by using a specific value of the parameter (irrespective of the data). Then, for this one parameter value, $\ell_n$ is zero identically, and in this one case we have the best overall model, but we will "pay" for this success at all other values of the parameter – if we are wrong then there may be a very heavy cost to using the empirical model $G_n$.

The technical framework used here is analogous to Phillips (1996) and Phillips and Ploberger (1996). In particular, we will maintain the following assumption among other conditions that will be detailed later.

**Assumption A0**

1. *The conditional probabilities $P_\theta(\cdot | \mathfrak{F}_{n-1})$ have densities $p_\theta(x_n | \mathfrak{F}_{n-1})$ (with respect to some dominating measure $\mu$ on $E$), our parameter space $\Theta \subset \mathbf{R}^k$, and the mapping $\theta \to \log p_\theta(x_n | \mathfrak{F}_{n-1})$ is twice continuously differentiable.*

2. *The score process component $\varepsilon_n(\theta) = \frac{\partial}{\partial \theta} \log p_\theta(x_n | \mathfrak{F}_{n-1})$ is square integrable. Define $B_n(\theta) = \sum_{1 \le i \le n} E_\theta(\varepsilon_i(\theta)\varepsilon_i(\theta)' | \mathfrak{F}_{i-1})$.*

3. *The prior distribution is proper with continuous density $\pi(\cdot)$ that is bounded away from the origin on every compact set $K$, so that $\inf_{\theta \in K} \pi(\theta) > 0$.*

The matrix $B_n$ in $\mathbf{A0}(2)$ is the conditional quadratic variation process of the score process $\sum_{i \leq n} \varepsilon_i(\theta)$. It can be regarded as one possible generalisation of the Fisher information matrix, a fact that we will more fully explore in following sections.

Phillips (1996) and Phillips and Ploberger (1996) show that if $P$ denotes the 'Bayesian' model (1) then, as $n \to \infty$, we have the asymptotic approximation

$$\log \frac{dP}{dP_\theta} \sim -\tfrac{1}{2} \log \det B_n(\theta) + (\theta - \hat{\theta}_n)' B_n(\theta)(\theta - \hat{\theta}_n)/2, \tag{3}$$

where $\hat{\theta}_n$ is the (normal) maximum-likelihood-estimator for $\theta$. What determines the order of magnitude of the terms on the right side of (3)? Clearly it is reasonable to assume that $\det B_n(\theta) \to \infty$, whereas the second summand is nothing else than the Wald–LM–LR test statistic for testing the parameter to be $\theta$. Asymptotic theory developed in recent years indicates that it is very plausible that — even under nonstationary circumstances — this statistic will remain $O_{P_\theta}(1)$. (For the case of general time series processes with some unit roots this is assured by the limit theorems in Park and Phillips, 1988, 1989). So, the term involving $\log \det B_n(\theta)$ in (3) will determine the order of magnitude of the loss that is due to the lack of information about the parameter. We will now show that it is only possible on a very small set of parameter values that, for arbitrary $\varepsilon > 0$, $\log \frac{dP}{dP_\theta} \geq -\frac{1-\varepsilon}{2} \log \det B_n(\theta)$ on a non-negligible event.

These two results have some interesting consequences for Bayesian models:

(i) Even from the point of view of our 'semi-classical' analysis, Bayesian models are *impossible to beat* from the predictive point of view.

(ii) The inevitable loss, $\log \det B_n(\theta)$, is easily seen to be dependent on the *dimension of the parameter space*. (In the stationary case, $B_n$ will asymptotically be of the form $B \cdot n$, therefore $\log \det B_n(\theta)$ will asymptotically be $\log \det(nB) = k \log n + O(1)$, where $k$ is the dimension of $B$). So, even the use of informative priors is no remedy against the curse of dimensionality. In other words it is essential to use parameters parsimoniously — a view that is commonly expressed by authors recommending methods for practitioners, e.g., Doan, Litterman and Sims (1984), West and Harrison (1989), Zellner and Min (1992).

In practice, it will often be the case that data will be explained not only in terms of their own past, but also by covariates. Under some reasonable assumptions, we can deal with this type of complication in our framework. Let us assume that our data $x_n$ consist of two 'components' as in $x_n = (y_n, z_n)$ (again, $y_n$ and $z_n$ can take values in arbitrary spaces). Suppose $y_n$ are the endogenous variables (i.e., the variables we want to explain) and $z_n$ are the exogenous variables, i.e. the variables we take as 'given'. Then, our 'models' will be conditional probability measures explaining $y_n$ by $z_n$, $x_{n-1}$, ..., $x_1$, since we do not want to model the exogenous variables. (In econometrics, such variables often reflect the outcome of governmental or political decisions and, while these decisions influence economic variables, it is usually not a feasible option to model these variables themselves (i.e., to make distributional assumptions about them).

The formalized concepts of exogenity discussed in Engle, Hendry and Richard (1983) have a long and successful tradition in econometrics and we are able to apply them here. The key step is a formalisation of the plausible assumption that the exogenous variables can be modeled without the exogenous ones. We therefore should have the following factorization

$$p_\theta(x_n|x^{n-1}) = q_\theta(y_n|z_n,\ x^{n-1})f(z_n|z^{n-1},\ x^{n-1}) \tag{4}$$

wherein the density factorizes into the (parametrized) conditional density of $y$ and the conditional density for $z$. Since the exogenous variables should be modeled without any reference to the model for the endogenous variables, their conditional density is not dependent on the parameters needed to describe the model for the endogenous ones.

Strictly speaking, we can think of (4) as a *definition* of exogenity (For a detailed discussion we refer to the article cited above). So, assuming we have given our parametrized family in terms of the conditional densities $q_\theta(y_n|z_n,\ x^{n-1})$, we can define $\mathfrak{F}'_{n-1} = \sigma(z_n,\ x^{n-1})$ and the 'models' as conditional probabilities from $\mathfrak{F}'_{n-1}$ to $y_n$. Then, we can think of constructing models $g(x_n|x^{n-1})$ for the whole process $x$ by modeling the conditional distribution of $y_n$ given $\{z_n,\ x^{n-1}\}$ and the conditional distribution of the $z_n$ component by its *true* density $f$. These models depend on the 'true' (and unknown) density for the exogenous variables, but it only influences them (and not the endogenous component $y$). Moreover, since we do not want to predict the $z$ component, the unknown character of the true density is of no importance to us. It is easily seen that, in the density ratios $dG/dP_\theta$, this — unknown — density cancels out. Therefore, we may, without a limitation in generality, assume that $\mathfrak{F}_{n-1} = \sigma\left(x^{n-1}\right)$, and, consequently, we are able to assume that our likelihoods are of the form

$$\log p_\theta(x_n, ..., x_1) = \sum_{i=1}^{n} \log q_\theta(y_n|z_n, x^{n-1}) \tag{5}$$

## 3   The Main Theorem

This section lays out our main result. Of central importance to our development will be an augmented space - * together with a $\sigma$-algebra $\mathfrak{F}^*$, which are defined as follows.

**Definition 1**   *Let - * $= \Theta \times$ - and let $\mathfrak{F}^*$ be the corresponding product $\sigma$-algebra of the Borel field of $\Theta$ and $\mathfrak{F}$. Analogously, let $\mathfrak{F}_n^*$ be the product $\sigma$-algebra of the Borel fields of $\Theta$ with $\mathfrak{F}_n$.*

This augmented space has some interesting properties. In particular, we can, for fixed $\theta \in \Theta$, extend our measures $P_\theta$ to - * by defining $P_\theta(A \times B) = I_A(\theta)P_\theta(B)$ for $A \subset \Theta$, $B \in \mathfrak{F}$ and then use standard measure theory to extend it to the whole product $\sigma$-algebra. (Here, $I_A(\cdot)$ is the indicator function of the set $A$.)

- $^*$ consists of pairs $(\theta, \omega)$, where $\theta \in \Theta$. We now consider the random variable (i.e., the mapping) attaching to each pair its first component, which we will denote for notational convenience by $\theta$, too. This random variable can be understood as the 'true' parameter, because the distribution of this random variable *under the measure* $P_\theta$ is trivial, viz., $P_\theta\left(\{(\theta, \omega) : \omega \in \text{-} \}\right) = 1$ and $P_\theta(\{(\vartheta, \omega) : \vartheta \in \Theta, \vartheta \neq \theta\}) = 0$. This concept of a "true" parameter, also makes sense for probability measures outside the set $\{P_\theta\}$.

We can, also extend the Bayesian mixture (1) to this probability space. Define for $A \subset \Theta$, $B \in \mathfrak{F}$ the measure $P(A \times B) = \int_A \pi(\theta) P_\theta(B) d\theta$ and then extend the measure to $\mathfrak{F}^*$. Restricting this measure to $\mathfrak{F}$ one easily sees that it is identical to (1). In what follows, we often need to do probability calculations with the measure $P$ (for example, we may need to show that certain random quantities are $O_P(1)$ as $n \to \infty$) and this formulation will then be very useful.

What is the advantage of this construction? Working with $\Theta \times \text{-}$ as the basic space enables us to consider the fundamental objects that we work with (e.g., likelihood processes) which are really continuous random fields indexed with $\theta$, as *simple random variables*. Indeed, a random field $Z_\theta$ (indexed by $\theta$) is just a family of measurable mappings from - into the real numbers. It is now an elementary task (if there exists a countable dense subset on $\Theta$) to show that the following statement holds. "For almost all $\omega \in \text{-}$, the mapping $\theta \to Z_\theta(\omega)$ is continuous" implies "the mapping $(\theta, \omega) \to Z_\theta(\omega)$ is (almost surely equal to) a measurable mapping". In the sequel, we will use this construction without further mentioning it. For most of the paper we will find this "random variable interpretation" of the likelihood process is better suited to our purposes. So we will, if not explicitly mentioned otherwise, assume that we are working on - $^*$ rather than - .

Besides Assumptions **A0,** our development relies on two 'high-level' assumptions, **A1** and **A2**, that are given below. **A1** simply guarantees that the information (in all parametric directions) contained in our experiment diverges to infinity when the sample size increases. This condition is the equivalent of a persistent excitation condition in regression models. Assumption **A2** postulates that, after we have "cut out" a small event, there exist measures that have a density of the order of magnitude of $1/(\det B_n(\theta))^{1/2}$. The existence of such a density in a very general class of econometric models is described in Phillips and Ploberger (1992, 1996).

**Assumption A1.** $\lambda_{\min}(B_n) \to \infty$ *a.s.* $(P_\theta)$ *for* $n \to \infty$, *where* $\lambda_{\min}(\cdot)$ *denotes the smallest eigenvalue.*

**Assumption A2.** *For every* $\eta > 0$ *there exist measures* $Q_n^{(\eta)}$ *on* $F_n$ *so that*

1. $\limsup_{n \to \infty} TV(Q_n^{(\eta)}, P^{(n)}) \leq \eta$, *where* $TV$ *denotes the total variation (or variational distance) between the measures.*

2. $\frac{dQ_n^{(\eta)}}{dP_\theta^{(n)}} \sqrt{\det B_n(\theta)} = O_P(1)$ *as* $n \to \infty$ *on a sequence of sets* $F_n \in \mathfrak{F}_n^*$ *for which* $\liminf_{n \to \infty} P(F_n) > 1 - \delta$ *for arbitrarily small* $\delta > 0$. *That is given* $\delta > 0$ *there exists*

$M_\delta$ such that

$$\liminf_{n \to \infty} P\left(F_n \cap \left[\frac{dQ_n^{(\eta)}}{dP_\theta^{(n)}}\sqrt{\det B_n(\theta)} < M\right]\right) > 1 - \delta,$$

where $P$ is our extended measure on $\mathfrak{F}_n^*$.

**Theorem 1** *Let Assumptions* A0, A1 *and* A2 *hold and let* $G$ *be an empirical "model measure". Then, for every compact* $K$, *the Lebesgue measure of*

$$\left\{\theta : P_\theta\left(\left[\log\left(\frac{dG^{(n)}}{dP_\theta^{(n)}}\right) \geq -\frac{1-\varepsilon}{2}\log\det B_n(\theta)\right]\right) \geq \alpha\right\} \cap K \qquad (6)$$

*converges to* 0 *as* $n \to \infty$.

**Proof of Theorem 1**  Let $G$ be a model measure. This model measure is defined on $(\text{-}, \mathfrak{F})$. But, it is easily seen that this measure can be extended to $(\text{-}^*, \mathfrak{F}^*)$ by defining $G(A \times B) = \int_A \pi(\theta)d\theta \cdot G(B)$ for $A \subset \Theta, B \subset \text{-}$. Analogously we can extend the $Q_n^{(\eta)}$ to $\mathfrak{F}_n^*$, too. To simplify notation, we just denote these extensions by $Q_n^{(\eta)}$ as well. Then $Q_n^{(\eta)}(A \times B) = \int_A \pi(\theta)d\theta \cdot Q_n^{(\eta)}(B)$ for $A \subset \Theta, B \subset \text{-}$.

We have to show that for all $\alpha, \varepsilon > 0$ and all compact $K$

$$\lambda\left(\left\{\theta \in K : P_\theta\left[\log\left(\frac{dG^{(n)}}{dP_\theta^{(n)}}\right) \geq -\frac{1-\varepsilon}{2}\log\det B_n(\theta)\right] \geq \alpha\right\}\right) \to 0,$$

where $\lambda(\cdot)$ is Lebesgue measure on $\Theta$.

Choose $\alpha, \varepsilon > 0$ and fix a compact $K$. Define the sets

$$C_n = \left\{\theta \in K : P_\theta\left[\log\left(\frac{dG^{(n)}}{dP_\theta^{(n)}}\right) \geq -\frac{1-\varepsilon}{2}\log\det B_n(\theta)\right] \geq \alpha\right\},$$

and

$$\Gamma_n = \left\{(\theta, \omega) : \theta \in C_n, \text{ and } \log\left(\frac{dG^{(n)}}{dP_\theta^{(n)}}(\omega)\right) \geq -\frac{1-\varepsilon}{2}\log\det B_n(\theta)(\omega)\right\}.$$

Then, with $\pi_0(K) = \inf_{\theta \in K} \pi(\theta) > 0$ we have $P(\Gamma_n) = \int_{C_n} P_\theta(\Gamma_n)\pi(\theta)d\theta \geq \alpha \cdot \pi_0(K) \cdot \lambda(C_n)$. Therefore, for the theorem to hold it is sufficient to show that $P(\Gamma_n) \to 0$. This assertion follows by showing, as we do below, that for an arbitrary $\eta > 0$ we have $\limsup_{n \to \infty} Q_n^{(\eta)}(\Gamma_n) \leq 7\eta$. Then, **A2**(1) gives the required result for $P(\Gamma_n)$.

First, Assumption **A2**(2) guarantees that $\frac{dQ_n^{(\eta)}}{dP_\theta^{(n)}}\sqrt{\det B_n(\theta)}$ remains $O_P(1)$. Therefore, there exists an $M_2 = M_2(\eta)$ for which, with $K_{1,n} = \left[\frac{dQ_n^{(\eta)}}{dP_\theta^{(n)}}\sqrt{\det B_n(\theta)} \leq M_2\right]$,

8

where $P(K_{1,n}) \geq 1 - \eta$. As $\limsup_{n\to\infty} TV(P^{(n)}, Q_n^{(\eta)}) \leq \eta$, there exists an $N_1 = N_1(\eta)$ such that for $n \geq N_1$, $\left| P(K_{1,n}) - Q_n^{(\eta)}(K_{1,n}) \right| < 2\eta$ and, consequently, $Q_n^{(\eta)}(K_{1,n}) \geq 1 - 3\eta$.

By Assumption **A1**, $\det B_n \to \infty$. Therefore, there exists an $N_2 = N_2(\eta)$ such that, with $K_{2,n} = \left[ (\det B_n)^{\varepsilon/2} \geq \frac{1}{\eta} M_2 \right]$, and $\varepsilon > 0$ arbitrary, $P(K_{2,n}) \geq 1 - \eta$. We can, without loss of generality, choose $N_2 \geq N_1$, and therefore $Q_n^{(\eta)}(K_{2,n}) \geq 1 - 3\eta$.

It is now sufficient to show that $\limsup_{n\to\infty} Q_n^{(\eta)}(\Gamma_n \cap K_{1,n} \cap K_{2,n}) \leq \eta$. Let $(\theta, \omega) \in \Gamma_n \cap K_{1,n} \cap K_{2,n}$ and let $n \geq \max(N_1, N_2)$. Since $(\theta, \omega) \in \Gamma_n$, we have $\frac{dG^{(n)}}{dP_\theta^{(n)}}(\omega) \geq \sqrt{(\det B_n(\theta))^{\varepsilon-1}}(\omega)$ and, since $\omega \in K_{1,n} \cap K_{2,n}$,

$$\frac{dP_\theta^{(n)}}{dQ_n^{(\eta)}}(\omega) = \frac{1}{\frac{dQ_n^{(\eta)}}{dP_\theta^{(n)}}(\omega)} \geq \frac{1}{M_2} \sqrt{\det B_n(\theta)}(\omega).$$

So, on $K_{1,n} \cap K_{2,n}$ we have

$$\frac{dG^{(n)}}{dQ_n^{(\eta)}}(\omega) = \frac{dG^{(n)}}{dP_\theta^{(n)}}(\omega) \cdot \frac{dP_\theta^{(n)}}{dQ_n^{(\eta)}}(\omega) \geq \frac{1}{M_2}(\det B_n)^{\varepsilon/2} \geq \frac{1}{\eta}.$$

Hence,

$$
\begin{aligned}
1 &\geq G(\Gamma_n \cap K_{1,n} \cap K_{2,n}) \\
&\geq \int_{\Gamma_n \cap K_{1,n} \cap K_{2,n}} \frac{dG^{(n)}}{dP_\theta^{(n)}} \cdot \frac{dP_\theta^{(n)}}{dQ_n^{(\eta)}} \cdot dQ_n^{(\eta)} \pi(\theta) d\theta \qquad (7) \\
&\geq \frac{Q_n^{(\eta)}(\Gamma_n \cap K_{1,n} \cap K_{2,n})}{\eta}.
\end{aligned}
$$

Setting $K_n = K_{1,n} \cap K_{2,n}$ and letting $K_n^c$ be the complement of $K_n$, we have

$$
\begin{aligned}
Q_n^{(\eta)}(\Gamma_n) &= Q_n^{(\eta)}(\Gamma_n \cap K_n) + Q_n^{(\eta)}(\Gamma_n \cap K_n^c) \\
&\leq Q_n^{(\eta)}(\Gamma_n \cap K_n) + Q_n^{(\eta)}(K_n^c) \\
&\leq \eta + 6\eta,
\end{aligned}
$$

which delivers the required result. $\blacksquare$

**Remark** In the inequality (7), the "$\geq$" must not be replaced by an "$=$", as it may be possible that $G$ is not absolutely continuous with respect to $P_\theta^{(n)}$, in which case $dG^{(n)}/dP_\theta^{(n)}$ is the absolutely continuous part of $G^{(n)}$ only.

**Discussion**   Theorem 1 is related to a result on minimal information loss in modelling that was proved by Rissanen (1986, 1987). Rissanen showed that if the generating mechanism for the data is a stationary process and some technical conditions are fulfilled, then the Lebesgue measure of the set

$$\left\{ \theta : -E_\theta \left( \log \frac{dG^{(n)}}{dP_\theta^{(n)}} \right) \leq \frac{1}{2} k \log n \right\} \tag{8}$$

converges to 0 for any choice of empirical model $G^{(n)}$. This theorem showed that whatever one's model, one can approximate (with respect to KL distance) the DGP of a stationary process no better *on average* than $\frac{1}{2} k \log n$. Thus, outside of a 'small' set of parameters we can get no closer to the truth than $\frac{1}{2} k \log n$ - the 'volume' of the set for which we can do better actually converges to zero.

Our result has a similar interpretation. Up to a 'small' exceptional set, the empirical model $G^{(n)}$ cannot come nearer to the true DGP than $\frac{1}{2} \log \det B_n$ as shown in (6). Since $G^{(n)}$ is arbitrary, the result tells us that there is a bound on how close any empirical model can come to the truth and that this bound depends on the data through $B_n$. It may well therefore be path dependent, rather than being reliant solely on the dimension of the parameter space as (8).

Not only is there a bound on how close we can come in empirical modelling to the true DGP, but the bound is attainable. Indeed, Phillips (1996) and Phillips and Ploberger (1996) show how to can construct empirical models for which

$$- \left( \log \frac{dG^{(n)}}{dP_\theta^{(n)}} \right) / (\log \det B_n) \rightarrow \frac{1}{2}. \tag{9}$$

These models can be formed by taking $G^{(n)}$ to be the Bayesian data measure $P^{(n)}$ for proper Bayesian priors. Or, in the case of improper priors, the models $G^{(n)}$ may be obtained by taking the conditional Bayes measures, given some preliminary set of $n_0$ observations (as in (2) above). The models can also be obtained by prequential methods, like those in Dawid (1984) and Phillips (1996).

Models that are 'better' than those which attain (9) must satisfy the inequality defined by the event

$$A_n = \left[ - \left( \log \frac{dG^{(n)}}{dP_\theta^{(n)}} \right) / (\log \det B_n) \leq \frac{1-\varepsilon}{2} \right] \tag{10}$$

for some $\varepsilon > 0$ at least somewhere in the probability space. However, if the probability of the event $A_n$ converges to zero, one cannot reasonably define $G^{(n)}$ to be better because the sample space over which the inequality (10) holds has negligible probability. Therefore, for a model to be essentially better, we must postulate the existence of an $\alpha > 0$ for which $P_\theta(A_n) \geq \alpha$, and then the probability of events such as $A_n$ is non negligible. What Theorem 1 tells us is that the set of such essentially better models has Lebesgue measure zero in the parameter space in $R^k$ as $n \rightarrow \infty$. In this well defined sense, we can generally expect to be able to do no better in modeling the DGP than to use the Bayesian models $P^{(n)}$.

10

# 4  Sufficient Conditions for Assumption A2

As it is stated, **A2** is a 'high-level' assumption. This section reformulates the assumption into more familiar terms and provides more primitive conditions for its validity. In earlier work (Phillips and Ploberger, 1996) the behaviour of the density of the Bayesian mixture measure (1) with respect to the true measure $P_\theta$ was investigated. It was shown there that, for a rather wide class of econometric models and under relatively weak regularity assumptions, the Bayesian data density $\frac{dP}{dP_\theta}$ is asymptotically proportional to $\pi(\theta) \frac{1}{\sqrt{\det B_n(\theta)}} \exp(\hat{\theta}_n - \theta)' B_n(\theta)(\hat{\theta}_n - \theta)/2)$ , where $\hat{\theta}_n$ is the maximum likelihood estimator for $\theta$. We now proceed to utilize these asymptotic results and some of the primitive conditions of that earlier work in validating **A2**. We start with the following assumption.

**Assumption B0**  $W_n(\theta) = (\hat{\theta}_n - \theta)' B_n(\theta)(\hat{\theta}_n - \theta) = O_{P_\theta}(1)$ *for Lebesgue almost all* $\theta \in \Theta$.

This assumption is plausible and can be expected to hold under quite general conditions. First, the statistic $W_n(\theta)$ is analogous to a Wald statistic and forms the basis of an asymptotic test that the parameter $\theta$ takes on a certain value. Under $P_\theta$, it is reasonable to suppose that $W_n(\theta) = O_{P_\theta}(1)$, although the critical values may well be nonstandard and, in some cases, even parameter dependent (this means dependent on $\theta$, here, as there are no extra nuisance parameters in our $P_\theta$). Obviously, the condition is fulfilled in the 'classical' case of stationary time series, but it has also been established in models with unit roots (Phillips and Durlauf, 1986) and with unit roots and cointegration (Park and Phillips, 1988, 1989). Note that one obvious implication of **B0** and the excitation condition **A1** is that $\hat{\theta}_n \to_p \theta$ $(P_\theta)$ for Lebesgue-almost all $\theta \in \Theta$. Thus, the MLE is consistent almost everywhere (Lebesgue measure) in the parameter space.

Together with Assumption **B0**, the results from Phillips (1996) give sufficient conditions (conditions C1-C7 in that paper) for **A2** to hold. They cover almost all 'classical' (i.e., asymptotically stationary) situations as well as cases with unit roots and cointegration. We will, however, go one step further. Here we are not so much interested in the data density itself, we only want to bound it from above. We can therefore use more convenient conditions to assure this. Central to our derivation is the assumption that the second derivative of the log likelihood function is continuous in a neighbourhood of $\theta$. Our main focus, in fact, is a small shrinking neighbourhood of $\theta$. In effect, the probability measures corresponding to parameters in this neighbourhood are contiguous to the original measure. In the 'classical' case, these neighbourhoods shrink with the order of $1/\sqrt{n}$.

**Assumption B1**  *The conditional log likelihood* $\log p_\theta(x_t|\mathfrak{F}_{t-1})$ *is twice continuously differentiable (in $\theta$) and* $\frac{\partial \varepsilon_{t,\theta}}{\partial \theta}$ *is integrable, where* $\varepsilon_t(\theta) = \partial \log p_\theta(x_t|\mathfrak{F}_{t-1})/\partial \theta$, *as before.*

11

Under **B1** and since $\frac{\partial \varepsilon_{t,\theta}}{\partial \theta}$ is integrable, we have $E_\theta \left( \frac{\partial \varepsilon_{t,\theta}}{\partial \theta} | \mathfrak{F}_{t-1} \right) + E_\theta(\varepsilon_{t,\theta} \varepsilon'_{t,\theta} | \mathfrak{F}_{t-1}) = 0$. Hence $\sum_{t \leq n} \frac{\partial \varepsilon_{t,\theta}}{\partial \theta} + B_n(\theta)$ is a $P_\theta$-martingale. As $B_n(\theta)$ increases monotonically and diverges (in view of **A1**), it is reasonable to assume that $\sum_{t \leq n} \frac{\partial \varepsilon_{t,\theta}}{\partial \theta} + B_n(\theta)$ is 'small' compared with $B_n(\theta)$, or, for each vector $h$, $\sum_{t \leq n} h' \frac{\partial \varepsilon_{t,\theta}}{\partial \theta} h + h' B_n(\theta) h = o(h' B_n(\theta) h)$. This requirement is a standard assumption in asymptotic theory, c.f. Hall and Heyde, 1980, Ch. 6.). In Phillips(1996) the requirement was assumed to hold uniformly in $h$, i.e.

$$\sup_{\|h\|=1} \left| \frac{\sum_{t \leq n} h' \frac{\partial \varepsilon_{t,\theta}}{\partial \theta} h + h' B_n(\theta) h}{h' B_n(\theta) h} \right| = o_{P_\theta}(1).$$

Denote by $\ell_t(\theta)$ the log likelihood function and by $\ell_t^{(1)}(\theta), \ell_t^{(2)}(\theta)$ its first and second $\theta$-derivatives. We reformulate the above requirement in the following form.

**Assumption B2**   *For Lebesgue-almost all $\theta \in \Theta$*

$$\sup_{\|h\|=1} \left| \frac{h' \ell_n^{(2)}(\theta) h + h' B_n(\theta) h}{h' B_n(\theta) h} \right| \to_{P_\theta} 0.$$

We also use another well-established asymptotic technique, namely the local approximation of the log-likelihood with a quadratic over 'shrinking' neighbourhoods (c.f. Phillips, 1996 and Kim, 1994). We have to be careful in making our assumptions about this phenomena, since we want to allow for generality and are especially interested in cases where the information matrix (i.e, $B_n(\theta)$) is neither asymptotically constant nor regular in the sense that its eigenvalues can have different orders of magnitude. To accomplish this, let $M > 0$ and define the following shrinking neighbourhood system of $\theta_0$

$$E_M(\theta_0) = \{\theta : (\theta_0 - \theta)' B_n(\theta)(\theta_0 - \theta) \leq M\}.$$

**Assumption B3**   *For all $M > 0$*

$$\sup_{\|h\|=1, \theta \epsilon E_M(\theta_0)} \left| \frac{h' \ell_n^{(2)}(\theta) h - h' \ell_n^{(2)}(\theta_0) h}{h' B_n(\theta_0) h} \right| \to_{P_{\theta_0}} 0.$$

Finally, we add the following technical requirement on the space $\Theta$.

**Assumption B4**   *The boundary of $\Theta$ (i.e., the difference between its closure and the interior) has Lebesgue-measure zero.*

Theorem 2 below gives sufficient conditions for **A2** in terms of these more primitive assumptions. Before stating the theorem, we give two technical lemmas that are useful in what follows. The first provides a formula for a restricted Radon Nikodym density in terms of mixture densities.

**Lemma RRN**   *Suppose we define for every set $F \in \mathfrak{F}_n^*$ the measure $\mu_F$ by $\mu_F(A) = P(A \cap F)$ and let $\nu_F$ be its restriction to $\mathfrak{F}_n$. Then*

$$\frac{d\nu_F}{dP_{\theta_0}^{(n)}} = \int_\Theta I_F(\theta) \frac{dP_\theta}{dP_{\theta_0}} \pi(\theta) d\theta. \tag{11}$$

**Proof**   Use a monotone class argument. Evidently, the lemma is valid for all sets

$$F = B \times C, B \subset \Theta, C \subset \text{-} . \tag{12}$$

Moreover, if it is true for sets $F'$, $F''$ with $F' \subset F''$, then it is valid for $F'' - F'$, too. Furthermore, if the relationship is true for a monotone increasing sequence of events $F_k, k = 1, 2, \ldots$, then it is true for its limit also. Therefore, the set of all sets $F$ for which the lemma is true is a Dynkin-system generated by the sets 12. As this generating set is $\cap$-stable, the Dynkin system is the whole $\sigma$-algebra, which proves the lemma.   ∎

The second lemma gives us a useful technique for converting $O_{P_\theta}$ bounds into $O_P$ bounds.

**Lemma P-BD**   *Suppose we are given two sequences of processes $E_n(\theta)$ and $F_n(\theta)$, for which $E_n(\theta) = O_{P_\theta}(E_n(\theta))$, for Lebesgue almost all $\theta \in \Theta$. Moreover, given $\varepsilon > 0$ and*

$$M(\varepsilon, \theta) < \infty \tag{13}$$

*for which*

$$P_\theta \left[ \left| \frac{E_n(\theta)}{F_n(\theta)} \right| \geq M(\varepsilon, \theta) \right] \leq \varepsilon, \tag{14}$$

*almost everywhere in $\theta$, it is further assumed that the bounding quantity $M(\varepsilon, \theta)$ is measurable in $\theta$. Then*

$$E_n(\theta) = O_P(E_n(\theta)) \tag{15}$$

*where $P = \int P_\theta \pi(\theta) d\theta$ is the mixture measure (1) and $\pi(\cdot)$ is a proper prior distribution on $\Theta$ with $\int \pi(\theta) d\theta = 1$.*

**Proof**   In view of (13) we can write $\Theta = \bigcup_{k \in \mathbf{N}} \{\theta : M(\varepsilon, \theta) < k\}$, at least up to a set of Lebesgue measure zero in $R^k$. Hence, by virtue of the integrability of $\pi(\cdot)$, we have

$$\lim_{k \to \infty} \int_{\{\theta : M(\varepsilon, \theta) \geq k\}} P_\theta \pi(\theta) d\theta = \lim_{k \to \infty} P[M(\varepsilon, \theta) \geq k] = 0.$$

For the last equation above, observe that $\theta$, and $M(\varepsilon, \theta)$ are random variables, the latter due to the measurability assumption on $M(\varepsilon, \theta)$.

It follows that for all $\varepsilon > 0$ we can find a $K(\varepsilon)$ so that

$$P[M(\varepsilon, \theta) \geq K(\varepsilon)] < \varepsilon. \tag{16}$$

To demonstrate (15) it is sufficient to show that for all $\varepsilon > 0$

$$P\left[\left|\frac{E_n(\theta)}{F_n(\theta)}\right| \geq K(\varepsilon)\right] \leq 2\varepsilon. \tag{17}$$

To show (17) holds, write

$$\left[\left|\frac{E_n(\theta)}{F_n(\theta)}\right| \geq K(\varepsilon)\right] \subseteq \left(\left[\left|\frac{E_n(\theta)}{F_n(\theta)}\right| \geq M(\varepsilon,\theta)\right] \cap [M(\varepsilon,\theta) < K(\varepsilon)]\right) \cup [M(\varepsilon,\theta) \geq K(\varepsilon)]. \tag{18}$$

Then, in view of the construction of $K(\varepsilon)$ in (16), the probability of the second event $[M(\varepsilon,\theta) \geq K(\varepsilon)]$ in (18) is $\leq \varepsilon$, whereas for the first event we have

$$
\begin{aligned}
&P\left(\left[\left|\frac{E_n(\theta)}{F_n(\theta)}\right| \geq M(\varepsilon,\theta)\right] \cap [M(\varepsilon,\theta) < K(\varepsilon)]\right) \\
&= \int P_\theta\left(\left[\left|\frac{E_n(\theta)}{F_n(\theta)}\right| \geq M(\varepsilon,\theta)\right] \cap [M(\varepsilon,\theta) < K(\varepsilon)]\right) \pi(\theta)\, d\theta \\
&= \int_{[M(\varepsilon,\theta)<K(\varepsilon)]} P_\theta\left(\left[\left|\frac{E_n(\theta)}{F_n(\theta)}\right| \geq M(\varepsilon,\theta)\right]\right) \pi(\theta)\, d\theta \\
&\leq \varepsilon \int \pi(\theta)\, d\theta = \varepsilon,
\end{aligned}
$$

where we use the fact that (14) holds for Lebesgue almost all $\theta$ by assumption. Summing these probabilities gives (17), and the result follows. ■

**Remark** The measurability assumption in Lemma P-BD seems quite mild and facilitates the conversion of $P_\theta$ probabilities of bounding events into $P$ probabilities. When we require this measurability assumption in future, we will simply say "with measurable bounds". An alternative approach would be to assume directly that the $O_{P_\theta}$ bounds hold uniformly in $\theta$, which is a more severe restriction and one that may be violated in some cases where limit distributions do not occur uniformly in the parameter space, as happens in some time series situations like those involving unit roots.

In the sequel and particularly in the proof of theorem 2 below, we often deal with inequalities between random variables defined on our augmented space $\Theta \times$ - which are not valid for all elements of $\Theta \times$ - . In such cases the following definition is useful.

**Definition 2** *Given random variables $X_1$, $X_2$ on $\Theta \times$ - , we say that*

$$X_1 \leq X_2 \text{ on a set } F$$

*if and only if*

$$I_F X_1 \leq I_F X_2.$$

14

We are now in a position to state our main result on sufficient 'primitive' conditions for **A2**.

**Theorem 2**  *Suppose Assumptions* **A0**–**A1** *and* **B0**–**B4** *are fulfilled with measurable bounds. Then, Assumption* **A2** *holds.*

**Proof**   We need to show that for every $\eta > 0$ we can approximate $P^{(n)}$ by measures $Q_n^{(\eta)}$ in such a way that **A2** holds, viz.,

1. $\limsup_{n\to\infty} TV(P^{(n)}, Q_n^{(\eta)}) \leq \eta$, and

2. $\frac{dQ_n^{(\eta)}}{dP_\theta^{(n)}}\sqrt{\det B_n(\theta)}$ remains $O_P(1)$ at least on a sequence of sets $F_n \in \mathfrak{F}_n^*$ for

which $\liminf_{n\to\infty} P(F_n) > 1 - \delta$ for arbitrarily small $\delta > 0$. That is given $\delta > 0$ there exists an $M_\delta$ for which

$$\lim\inf_{n\to\infty} P\left(F_n \cap \left[\frac{dQ_n^{(\eta)}}{dP_\theta^{(n)}}\sqrt{\det B_n(\theta)} < M_\delta\right]\right) > 1 - \delta.$$

Choose $\eta > 0$. Then, in view of **B0** we can find $M = M(\eta)$ so that

$$\lim\inf_{n\to\infty} P([\hat\theta_n \in E_M(\theta)]) \geq 1 - \eta.$$

Define the events $F_n^{(i)} \in \mathfrak{F}_n^*$, $i = 1, 2, 3, 4$, as follows:

$$F_n^{(1)} = [\hat\theta_n \in E_M(\theta)] \cap [E_{2M}(\theta) \subset \Theta], \qquad (19)$$

$$F_n^{(2)} = [-(\theta - \hat\theta_n)'\ell_n^{(2)}(\hat\theta_n)(\theta - \hat\theta_n) \leq 4M], \qquad (20)$$

$$F_n^{(3)} = \left[\sup_{\|h\|=1, \vartheta \in E_M(\theta)} \left|\frac{h'\ell_n^{(2)}(\vartheta)h - h'\ell_n^{(2)}(\theta)h}{h'B_n(\theta)h}\right| < \frac{1}{16}\right], \qquad (21)$$

$$F_n^{(4)} = \left[\sup_{\|h\|=1} \left|\frac{h'\ell_n^{(2)}(\theta)h + h'B_n(\theta)h}{h'B_n(\theta)h}\right| < \frac{1}{16}\right], \qquad (22)$$

and then set $F_n = F_n^{(1)} \cap F_n^{(2)} \cap F_n^{(3)} \cap F_n^{(4)}$. It is apparent that that $F_n \in \mathfrak{F}_n^*$ (and the same applies for $F_n^{(i)}$, $i = 1, 2, 3, 4$). It is important to understand that these sets are all subsets of $\Theta \times$ - .

Assumptions **B2** and **B3** imply that $\lim_{n\to\theta} P(F_n^{(3)} \cap F_n^{(4)}) = 1$. From the defining properties of the $F_n^{(i)}$ and $E_M(\theta)$ it can easily be seen that $F_n^{(1)} \cap F_n^{(3)} \cap F_n^{(4)} \subset F_n^{(2)} \cap F_n^{(3)} \cap F_n^{(4)}$. Therefore,

$$\lim\inf_{n\to\infty} P(F_n) \geq \lim\inf_{n\to\infty} P(F_n^{(1)} \cap F_n^{(3)} \cap F_n^{(4)}) \geq \lim\inf_{n\to\infty} P(F_n^{(1)}) \geq 1 - \eta.$$

Assumption B4 guarantees that $\lim_{n\to\infty} P[E_{2M}(\theta) \subset \Theta] = 1$.

Now define the measure $R_n^{(\eta)}$ on $\mathfrak{F}_n^*$ by $R_n^{(\eta)}(A) = P(F_n \cap A)$ and let $Q_n^{(\eta)}$ be its *restriction* on $\mathfrak{F}_n$. Then $TV(P^{(n)}, Q_n^{(\eta)}) \leq TV(P|\mathfrak{F}_n^*, R_n^{(\eta)}) = 1 - P(F_n)$, which shows that the first requirement of Assumption **A2** is fulfilled.

For the second part of **A2**, we have to compute $\frac{dQ_n^{(\eta)}}{dP_\theta^{(n)}}$. In the proof that follows, we will use the fact that the $Q_n^{(\eta)}$ are restrictions of the measures $R_n^{(\eta)}$. For all $A \in \mathfrak{F}_n$ we have $Q_n^{(\eta)}(A) = R_n^{(\eta)}(A) = P(A \cap F_n) = \int P_\theta(A \cap F_n)\pi(\theta)d\theta$. From this representation, the density can be computed easily by using (11) from Lemma RRN. In particular, for a given $\theta \in \Theta$, we have

$$\frac{dQ_n^{(\eta)}}{dP_\theta^{(n)}} = \int_\Theta I_{F_n}(\vartheta, .)\frac{dP_\vartheta^{(n)}}{dP_\theta^{(n)}}\pi(\vartheta)\,d\vartheta.$$

We now need to show that for $(\theta, \cdot)$ on $F_n$

$$\sqrt{\det B_n(\theta)} \int_\Theta I_{F_n}(\vartheta, \cdot)\frac{p_n(\vartheta)}{p_n(\theta)}\pi(\vartheta)\,d\vartheta = O_P(1), \quad \text{as } n \to \infty \tag{23}$$

where $p_n(\vartheta) = dP_\vartheta^{(n)}/d\mu$ is the density of $P_\vartheta^{(n)}$ and, similarly, $p_n(\theta)$. For $I_{F_n}(\vartheta, \cdot)$ to be nonzero, it follows from the construction of the set $F_n = F_n^{(1)} \cap F_n^{(2)} \cap F_n^{(3)} \cap F_n^{(4)}$ that

$$-(\vartheta - \hat\theta_n)'\ell_n^{(2)}(\hat\theta_n)(\vartheta - \hat\theta_n) \leq 4M,$$

which allows us to restrict the domain of integration accordingly.

It is easily seen from the definitions of $F_n^{(3)} \cap F_n^{(4)}$ and $F_n^{(1)}$ that on $F_n$

$$B_n(\theta) \leq 4(-\ell_n^{(2)}(\hat\theta_n)) \tag{24}$$

(in the usual partial ordering of non negative definite matrices), so that

$$\det B_n(\theta) \leq \det(4(-\ell_n^{(2)}(\hat\theta_n))). \tag{25}$$

Both (24) and (25) should be understood as inequalities between random variables defined on $\Theta \times$- . Thus, (24) means that if $(\omega, \theta) \in F_n$ then $B_n(\theta)(\omega) \leq 4(-\ell_n^{(2)}(\hat\theta_n))(\omega)$.

Moreover, we can use (19)–(22) to derive approximations for the second derivative of the log-likelihood. In particular, on $F_n$ we have

$$\sup_{\|h\|=1, \vartheta \in E_{4M}(\theta)} \left|\frac{h'\ell_n^{(2)}(\vartheta)h - h'\ell_n^{(2)}(\theta)h}{h'B_n(\theta)h}\right| \leq \frac{1}{16},$$

and, as $\hat\theta_n \in E_M(\theta)$, we also have

$$\sup_{\|h\|=1, \vartheta \in E_M(\theta)} \in \left|\frac{h'\ell_n^{(2)}(\vartheta)h - h'\ell_n^{(2)}(\hat\theta_n)h}{h'B_n(\theta)h}\right| \leq \frac{1}{8},$$

16

and, therefore, (using(24)) on $F_n$

$$\sup_{\|h\|=1,\vartheta\in E_M(\theta)}\left|\frac{h'\ell_n^{(2)}(\vartheta)h - h'\ell_n^{(2)}(\hat{\theta}_n)h}{h'\ell_n^{(2)}(\hat{\theta}_n)h}\right| \leq \frac{1}{2}.$$

We may conclude that for $\vartheta\in E_M(\theta)$, and all vectors $h$ we have on $F_n$

$$\frac{1}{2}h'\ell_n^{(2)}(\hat{\theta}_n)h \leq h'\ell_n^{(2)}(\vartheta)h \leq \frac{3}{2}h'\ell_n^{(2)}(\hat{\theta}_n)h.$$

As $E_M(\theta)$ is convex, we can use the Taylor expansion to conclude that for $\vartheta\in E_M(\theta)$ on $F_n$

$$\ell_n(\vartheta) \leq \ell_n(\hat{\theta}_n) + \frac{1}{4}(\vartheta-\theta)'\ell_n^{(2)}(\hat{\theta}_n)(\vartheta-\theta),$$

and

$$\ell_n(\hat{\theta}_n) \leq \ell_n(\theta) - \frac{3}{4}(\hat{\theta}_n-\theta)'\ell_n^{(2)}(\hat{\theta}_n)(\hat{\theta}_n-\theta).$$

As $\frac{dP_\vartheta^{(n)}}{dP_\theta^{(n)}} = \exp(\ell_n(\vartheta)-\ell_n(\theta))$, we therefore have the following inequality on $F_n$

$$\frac{dP_\vartheta^{(n)}}{dP_\theta^{(n)}} \leq \exp\left(\tfrac{1}{4}(\vartheta-\theta)'\ell_n^{(2)}(\hat{\theta}_n)(\vartheta-\theta)\right)\exp\left(-\tfrac{3}{4}(\hat{\theta}_n-\theta)'\ell_n^{(2)}(\hat{\theta}_n)(\hat{\theta}_n-\theta)\right).$$

Let $\pi_n = \sup_{\pi\in E_{4M}}\pi(\theta)$. Then

$$\int \frac{dP_\vartheta^{(n)}}{dP_\theta^{(n)}}\pi(\vartheta)d\vartheta$$

$$\leq \exp\left(-\tfrac{3}{4}(\hat{\theta}_n-\theta)'\ell_n^{(2)}(\hat{\theta}_n)(\hat{\theta}_n-\theta)\right)\int\exp\left(\tfrac{1}{4}(\vartheta-\theta)'\ell_n^{(2)}(\hat{\theta}_n)(\vartheta-\theta)\right)d\vartheta\pi_n \quad (26)$$

The first factor in (26) is $O_{P_\theta}(1)$ for Lebesgue-almost all $\theta$ due to assumptions **B0** and **B2.** It follows from Lemma P-BD and the measurability of the bound that this first factor on the right side of (26) is also $O_P(1)$ as $n\to\infty$. The second factor of (26) equals $C/\sqrt{\det(-\ell_n^{(2)}(\hat{\theta}_n))}$, where $C$ is a universal normalizing factor depending only on the dimension of $\theta$. Inequality (25) shows that, on $F_n$, $\det(-\ell_n^{(2)}(\hat{\theta}_n)) \geq$ Const $\cdot \det B_n(\theta)$ , which proves (23) and then **A2**(2) is established. ∎

## 5   Gaussian Models

This section deals with the important practical example of conditional Gaussian models. Under rather general conditions, we will show that these models satisfy Assumption **A2**. In particular, we do not limit ourselves to cases where the limiting distribution of the MLE is a mixture of Gaussian processes. For the theory to be useful in econometric applications that include unit roots and cointegration, one has

to include models where the limiting decision problem involves diffusion processes. To permit extensions to such situations, we do require some functional limit theory to be fulfilled. But the conditions are relatively mild and, as shown in Park and Phillips (1988, 1989), they are fulfilled for all models of practical interest.

The model class to be considered is prescribed by the systems equation

$$y_t = \Pi(\beta)z_t + u_t \tag{27}$$

where $y_t$ is a $k$-vector of *endogenous* variables, $z_t$ is a $K$-vector of *exogenous* or *predetermined* (i.e., $\mathfrak{F}_{t-1}$-measurable) variables, $\beta$ is a parameter vector, and $u_t =_d$iid $N(0, \Sigma)$ where $\Sigma = \Sigma(\gamma)$, i.e., we allow $\Sigma$ to depend on a parameter vector $\gamma$ that is to be estimated.

Let us now assume the following:

**Assumption C1**  *The parameter space $\Theta = \{(\beta, \gamma) : \beta \in \Theta_1, \gamma \in \Theta_2\}$ with $\Theta_1 \subset \mathcal{R}^\ell$, $\Theta_2 \subset \mathcal{R}^p$ and both sets are open and their boundaries have Lebesgue measure zero. Furthermore, the functions $\beta \to \Pi(\beta)$ and $\gamma \to \Sigma(\gamma)$ are twice continuously differentiable. Moreover, $\Sigma(\gamma)$ is (for Lebesgue-almost all $\gamma$) nonsingular.*

**Assumption C2**  *Both parameters are locally identified, i.e. the first derivatives of $\Pi$ and $\Sigma$ (with respect to $\beta$ and $\gamma$) have maximal rank (i.e., $\ell$ and $p$, respectively).*

**Assumption C3**  *For Lebesgue almost all $\theta$, there exist orthogonal matrices $O_n = O_n(\theta)$ and diagonal matrices $D_n(\theta) = D_n = \operatorname{diag}(\lambda_{i,n})$ such that $\liminf_{i,n} \lambda_{i,n} > 0$, and the random variables $W_n = \frac{1}{\sqrt{n}} \sum_{t \le n} D_n^{-1} O_n' z_t u_t'$ and $A_n = \frac{1}{n} \sum D_n^{-1} O_n' z_t z_t' O_n D_n^{-1}$ converge jointly in distribution. In particular, $(W_n, A_n) \to_d (W, A)$, where $W$ and $A$ are random elements and $A$ is positive definite (almost surely $P_\theta$).*

Then we have the following result which validates our main theorem under these conditions.

**Theorem 3**  *If the model (27) satisfies assumptions **C1**–**C3,** and all $O_{P_\theta}$ bounds are measurable in $\theta$, then Assumption **A2** holds.*

The proof of theorem 3 is lengthy and involves several technical lemmas. It is therefore given in the Appendix (see section 8.2).

# 6  Forecasting with Structural Linear Models

In this section we apply the above results to derive bounds for the quality of the prediction in linear models. In particular, we seek to determine how close to the optimal predictor we can come using empirical models, i.e. models in which the parameters have been estimated.

We consider a standard linear econometric model of the form

$$\Gamma y_t = B x_t + u_t \tag{28}$$

where $B$ and $\Gamma$ are the (partially unknown) parameter matrices, the $k$-vector $y_t$ contains *endogenous* variables and the $h$-vector $x_t$ consists of *exogenous* and *predetermined* (i.e., $\mathfrak{F}_{t-1}$-measurable) variables. So, $\Gamma$ is a $k \times k$-matrix, and $B$ is a $k \times h$-matrix. The set up includes traditional simultaneous equation models as well as VAR models.

Let us assume that $u_t$ are i.i.d $N(0, \Sigma)^1$ and independent of $x_t$. The conditional distribution of $y_t$ given $\mathfrak{F}_{t-1}$ will be denoted by the measure $G(\Gamma^{-1} B x_t, \Sigma)$. Then, if all the parameters were known, the best prediction for $y_t$ would be

$$\widetilde{y}_t = \Gamma^{-1} B x_t \tag{29}$$

and the *unavoidable* error $y_t - \widetilde{y}_t = \Gamma^{-1} u_t$ is distributed $N(0, \Gamma^{-1} \Sigma^{-1} \Gamma^{-1})$. In general, however, one has to estimate the matrices $B$ and $\Gamma$. Therefore, it is not possible to compute $\widetilde{y}_t$ and, in practice, one has to use another predictor for $y_t$ — say $\widehat{y}_t$ (generated, for instance, by plugging in estimates of $B$ and $\Gamma$ in (29). For our analysis, we do not have to be concerned with how this prediction is constructed, as long as it is $\mathfrak{F}_{t-1}$-measurable.

Our object is to investigate the asymptotic behaviour of the weighted forecast mean square difference

$$\Delta_n = \sum_{t=n_0}^{n} \{(y_t - \widetilde{y}_t)' \Sigma^{-1} (y_t - \widetilde{y}_t) - (y_t - \widehat{y}_t)' \Sigma^{-1} (y_t - \widehat{y}_t)\},$$

where $n_0$ is some point of initialization of the forecasts and where to simplify notation in what follows we can set $n_0 = 1$, with no loss of generality. In particular, we will show that there exists a number $K$ (depending on the degree of nonstationaritiy and the number of cointegrating relationships in $x_t$) which has the property that for Lebesgue almost all parameters and for all $\varepsilon > 0$

$$P\left(\left[\Delta_n \geq -\frac{1-\varepsilon}{2} K \log n\right]\right) \to 0.$$

This result shows the inherent advantage of the approach we are taking. Our generalization of Rissanen's theorem enables us to cover the case of prediction errors when the regressors are nonstationary. Interestingly, as we will see, in these cases something new happens. The additional errors do not follow the classical (number of parameters)*(logarithm of sample size) rule. Instead, in our new rule, we have

---

[1]This distributional assumption may seem to be restrictrive. However, we want to derive lower bounds for the prediction error due to the fact that we have to estimate parameters. In general, one does maintain specific asumptions about the distribution of the $u_t$ to obtain an optimal predictor. Our bounds are valid for all situations where Gaussian errors are not excluded-

to multiply the number of parameters by an additional factor that is essentially determined by the number and type of the trends in the regressors.

Before formulating the prediction theorem we make our assumptions specific. We assume that we have given a model of the form (28) and that the parameters are certain coefficients of $B$ and $\Gamma$, with the remaining coefficients being known by way of normalization and identifying restrictions.

**Assumption D1**  *The parameter space is given by the elements $B_{i,j}$, $(i,j) \in M_1$ and $\Gamma_{i,j}$, $(i,j) \in M_2$. All the other coordinates are known. Moreover, we assume that $M_1$ and $M_2$ are such that all of the identification assumptions of the preceding section are fulfilled.*

The problem we are dealing with is just another formulation of the usual identification problem for structural models. In the notation of the previous section $\Pi = \Gamma^{-1}B$ and therefore
$$d\Pi = -\Gamma^{-1}d\Gamma\Pi + \Gamma^{-1}dB \tag{30}$$
For our identification condition **C2** to be fulfilled for Lebesgue almost all parameters it is well known that the following necessary and sufficient conditions must be true.

1. $\Gamma$ is nonsingular for almost all parameters

2. For each $i$ such that $1 \le i \le k$ define index sets corresponding to the included variables (or coefficients) as follows:

$$M_1(i) = \{j : (i,j) \in M_1\}, \tag{31}$$

and

$$M_2(i) = \{j : (i,j) \in M_2\}. \tag{32}$$

Then, for Lebesgue almost all parameters the following rank-condition holds:

For all $i$ such that $1 \le i \le k$ the set of $h-$vectors $\hspace{2em}$ (33)
$$\{e_j : j \in M_1(i)\} \cup \{\pi_j : j \in M_2(i)\} \text{ are linearly independent,}$$

where the $e_j$ are $h-$vectors with all components zero except the $j$-th component which is unity, and $\pi_j$ is the $j$-th row of $\Pi$.

**Assumption D2**  *Any linear combination of the components of $x_t$ is either* stationary and ergodic *or* integrated of order one.[2]

Further, we define for all $a \in \mathbf{R}^h$ the process $e_t(a)$ to be either $a'x_t$ — if $a'x_t$ is stationary — or $\Delta a'x_t = a'x_t - a'x_{t-1}$ — if $a'x_t$ is nonstationary. Then the process

---

[2]Following convention, a process is said to be integrated of order one, or $I(1)$, if its first difference is stationary and has non zero spectral density at the origin. The first difference is in this event said to be $I(0)$.

$e_t(a)$ is stationary in both cases[3]. We can therefore (if we assume that the processes are purely nondeterministic) apply Wold's decomposition theorem and conclude that

$$e_t(a) = \sum_{i=0}^{\infty} c_i u_{t-i} = c(L)u_t, \tag{34}$$

where $u_{t-i}$ is white noise with variance $\sigma_a^2$. Clearly, the constants $c_i$ as well as the $u_t$ depend on $a$. Nevertheless, we can make the following assumption:

**Assumption D3** *For every $a \in \mathbf{R}^h$ the process $e_t(a)$ either is constant or in its Wold-decomposition (34) the following holds true:*

$$\sum_{i=0}^{\infty} i^{\frac{1}{2}}|c_i| < \infty, \tag{35}$$

*and*

$$c(1) = \sum_{i=0}^{\infty} c_i \neq 0. \tag{36}$$

Assumption **D3** guarantees that the autocorrelations between the components of $e_t$ converge to zero fast enough to assure the continuity of the spectral density of $e_t$. Further, for $a \neq 0$, $e_t(a) \neq 0$ and partial sums of the $e_t(a)$ may be assumed to satisfy a functional central limit theorem. That is, as a function of $z$, with $0 \leq z \leq 1$, we have

$$\frac{1}{\sqrt{n}} \sum_{t=1}^{[nz]} e_t(a) \rightarrow_d c(1)\sigma_a W(z), \tag{37}$$

where $W(z)$, $0 \leq z \leq 1$ is a standard Wiener process. The functional law (37) is known to hold under under (35) under quite weak conditions on $u_t$ (see Phillips and Solo, 1992).

Moreover, it is easily seen that (36) guarantees the strict positivity of the long term variance (i.e., $c(1)^2 \sigma_a^2 > 0$ ) and this implies that the variance of the nonstationary linear combinations increases linearly with time. For this study, we restrict ourselves to 'genuine' $I(1)$ processes and exclude processes that may be fractionally integrated.

For the formulation of Theorem 4 below we need to introduce another concept, which we call the "total degree of integration". While it is easy to clasify scalar processes as $I(1)$ or $I(0)$, this classification is, for our purposes, too crude in the multivariate case, where there may be some unit roots but not necessarily a full set. Heuristically, it seems reasonable to think of a bivariate process (say) with two independent integrated processes as being 'more' integrated than a bivariate process with one integrated component and one stationary component. It turns out that this concept of degree of integration plays a major role in determining empirical limits on forecasting ability.

---

[3]The function $a \rightarrow e_t(a)$ is discontinuous in some cases (e.g., if there are cointegrating relationships present in the original process).

**Definition 3** *Let $z_t$ be a vector process satisfying Assumptions **D1**–**D3**. Assume $z_t$ has $n_{\text{stat}}$ stationary components[4], $n_{\text{coint}}$ cointegrating relationships and $m$ components in total. Then, the 'total degree of integration' of $z_t$, (written as $TI(z_t)$) is defined as follows. If no component of $z_t$ contains a deterministic trend, then*

$$TI(z_t) = n_{\text{stat}} + n_{\text{coint}} + 2(m - (n_{\text{stat}} + n_{\text{coint}})).$$

*If at least one of the components of $z_t$ contains a linear trend then*

$$TI(z_t) = n_{\text{stat}} + n_{\text{coint}} + 2(m - (n_{\text{stat}} + n_{\text{coint}})) + 3.$$

*If at least one of the components of $z_t$ contains a time polynomial of degree $p$ then*

$$TI(z_t) = n_{\text{stat}} + n_{\text{coint}} + 2(m - (n_{\text{stat}} + n_{\text{coint}})) + (2p + 1).$$

To the extent that $TI(z_t)$ differs from $n_{\text{stat}} + n_{\text{coint}}$, it measures how many linear independent integrated components (or stochastic trends) are present in $z_t$. In addition, as we will see, $TI(\cdot)$ also determines the order of growth of the information matrix associated with the variables $z_t$.

It will be convenient in our following development to introduce some new notation to enable us to work equation by equation. In particular, we define for each $i$ with $1 \le i \le k$ new processes $r_t^{(i)}$. The dimension of $r_t^{(i)}$ is set as the sum of the number of elements of $M_1(i)$ and $M_2(i)$ (defined in (31) and (32)). Then, we define for each element of $M_1(i)$ and each element of $M_2(i)$ a component of $r_t^{(i)}$ as follows: for $j \in M_1(i)$ we let the component equal $(x_t)_j$, the $j$-th component of $x_t$, and for $j \in M_2(i)$ we let the component be $(\Pi x_t)_j$, the $j$-th component of the vector $\Pi x_t$.

Heuristically, this construction can be described in the following way. We consider all the parameters to be estimated in equation $i$. For each parameter in $B$ we take the corresponding random variable as a component, and for each parameter in $\Gamma$ we take the corresponding component from the reduced form.

With this definition in hand, we can formulate our theorem on feasible empirical limits to forecasting and proximity to the optimal predictor when there are parameters to be estimated. Since the proof is lengthy it is given in the Appendix (see section 8.3).

**Theorem 4** *Suppose we have given a model* (28) *satisfying Assumptions **D1**-**D3**. Fix $\Sigma = E(u_t u_t')$. Then, for all strictly positive $\alpha$ and $\varepsilon$ the Lebesgue measure of the set of parameters*

$$\left[\theta = \{(B_{i,j}, \Gamma_{k,h})_{(i,j) \in M_1 \text{ and } (k,h) \in M_2}\} \text{ such that } P_\theta\left[\Delta_n \ge -\frac{1 - \varepsilon}{2} K \log n\right] \ge \alpha\right]$$

*converges to zero, where $K = \sum_i TI(r_t^{(i)})$ (and the $r_t^{(i)}$ are defined above).*

---

[4]*For convenience we also allow for constant components. The nonsingularity condition in Assumption **D2**, however, restricts us to just one possible constant component.*

**Remark** In the case of univariate $y_t$ the assumption on $\Sigma$ is harmless. It is an easy consequence of the results from Phillips and Ploberger (1994) that the bound is sharp, since the MLE of the coefficients does not depend on $\Sigma$. In the multivariate case, we usually have no information about $\Sigma$ and it will generally affect the MLE. However, this fact does not make our bound any less valid. For, even if one knew $\Sigma$, it would be impossible to get a better forecast! It does remain to show that this bound is attainable - for some special cases, see Gerencser and Rissanen (1992). The general nonstationary case is, to the best of our present knowledge, still an open problem that is of obvious interest and importance. We are optimistic that there will be a positive solution of the problem.

## 7 Conclusion

Theorem 1 and Rissanen's result (8) justify a certain amount of skepticism about models with large numbers of parameters. In the stationary case, it is relatively easy to compare the 'loss' from parameter estimation in different parameter spaces. According to Rissanen's result, the loss due to parameter estimation is essentially determined by the dimension of the parameter space. In this case, the *minimum achievable distance* of an empirical model to the DGP increases linearly with the number of parameters. In the presence of nonstationarities, however, the situation changes. It is not the dimension of the parameter space that determines the distance of the model to the true DGP, but the order of magnitude of the sample Fisher information matrix. All the commonly arising cases lead to asymptotic expressions of the form

$$\log \det B_n \sim \left( \sum_{i=1}^{k} \alpha_i \right) \log n \tag{38}$$

for the sample information and $\alpha_i \geqq 1$ with inequality occuring for at least one element $i$. In particular, $\alpha_i = 2$ for stochastic trends and $\alpha_i = 3$ for a linear deterministic trend. In such cases, the distance of the empirical model to the DGP increases *faster* than in the traditional case. In effect, when nonstationary regressors are present, it appears to be even more important to keep the model as simple as possible. An additional stochastic trend in a linear regression model will be twice as expensive as a stationary regressor in terms of the marginal increase in the nearest possible distance to the DGP and a linear trend three times more expensive. Although nonstationary regressors embody a powerful signal and have estimated coefficients that display faster rates of convergence than those of stationary regressors, they can also be powerfully wrong in prediction when inappropriate and so the loss from including nonstationary regressors is correspondingly higher. In a very real sense, therefore, the true DGP turns out to be more elusive when there is nonstationarity in the data.

The above remarks apply irrespective of the modelling methodology that is involved. Neither Bayesian nor classical techniques can overcome the bound on empirical modelling. The bound can be improved only in 'special' situations, like those

where we have extra information about the true DGP and do not have to estimate all the parameters. For instance, we may 'know' that there is a unit root in the model, or by divine inspiration we may hit upon the right value of a parameter and decide not to estimate it.

As we have seen, these results that delimit the achievable proximity to the true DGP in empirical modelling have counterparts in terms of the capacity of empirical models to capture the good properties of the optimal predictor (i.e. the predictor that uses knowledge of the DGP and, in particular, the values of its parameters). Increasing the dimension of the parameter space carries a price in terms of the quantitative bound of how close we can come to replicating the optimal predictor. Furthermore, this price goes up when we have trending data and when we use trending regressors.

# 8 Appendix

## 8.1 Preliminary Development and Lemmas

For the analysis in Section 5 it will be convenient to define $H(\gamma) = \Sigma(\gamma)^{-1}$, and then the log likelihood function for model (27) can be expressed as

$$\ell_n(\beta, \gamma) = \frac{n}{2} \log \det H(\gamma) - \frac{1}{2} \sum_{t \leq n} (y_t - \Pi(\beta)z_t)' H(\gamma)(y_t - \Pi(\beta)z_t).$$

Some elementary calculations yield the following results about the conditional quadratic variation matrix $B_n$ in this case.

**Lemma A1**

1. $B_n$ is block-diagonal: $B_n = \begin{pmatrix} B_n^{(\beta)} & 0 \\ 0 & B_n^{(\gamma)} \end{pmatrix}$;

2. $\lim_{n \to \infty} \frac{1}{n} B_n^{(\gamma)}$ is constant, nonsingular and a continuous function of $H(\gamma) = \Sigma(\gamma)^{-1}$.

3. $(B^{(\beta)})_{i,j} = \operatorname{tr}\left(\sum_{t \leq n} z_t z_t' \cdot \frac{\partial \Pi}{\partial \theta_i}' H \frac{\partial \Pi}{\partial \theta_j}\right).$

In the sequel, we often use bounds for matrix products of the form $\operatorname{tr}(AB)$ and the following result, whose proof is straightforward, is a useful tool.

**Lemma A2** *Let $A$, $B$, $C$ be nonnegative definite matrices with $B \leq C$. Then $\operatorname{tr}(AB) \leq \operatorname{tr}(AC)$.*

Using the notation of **C3**, define the matrices $S_n$ and $\Psi_n$ by $S_n = O_n(D_n \cdot D_n)O_n'$ and $(\Psi_n)_{i,j} = n \cdot tr(S_n \frac{\partial \Pi}{\partial \theta_i}' \frac{\partial \Pi}{\partial \theta_j}).$

24

**Lemma A3** *For every $\eta > 0$ there exist $a(\eta)$, $A(\eta) > 0$ for which*

$$P\left[a(\eta)\Psi_n \le B_n^{(\beta)} \le A(\eta)\Psi_n\right] \ge 1 - \eta.$$

**Proof** From Assumption **C3** ,we may conclude that for every $\eta > 0$ there exist $c(\eta)$, $C(\eta) > 0$, such that with $K_n = \left[c(\eta)I \le \frac{1}{n}\sum D_n^{-1}O_n'z_tz_t'O_nD_n^{-1} \le C(\eta)I\right]$, we have $P(K_n) \ge 1 - \eta/2$. Then, on $K_n$, $c(\eta)S_n \le \frac{1}{n}\sum z_tz_t' \le C(\eta)S_n$.

Let $h \in R^m$. Then, from lemma A1, $h'B_nh = \text{tr}\left((\sum_t z_tz_t')(\sum_{i,j}h_ih_j\frac{\partial\Pi}{\partial\theta_i}'H\frac{\partial\Pi}{\partial\theta_j})\right)$, and lemma A2 shows that on $K_n$

$$c(\eta)\cdot\text{tr}\left[S_n\left(\Sigma_ih_ih_j\frac{\partial\Pi}{\partial\theta_i}'H\frac{\partial\Pi}{\partial\theta_j}\right)\right] \le h'B_nh \le C(\eta)\cdot\text{tr}\left[S_n\left(\Sigma_{i,j}h_ih_j\frac{\partial\Pi}{\partial\theta_i}'H\frac{\partial\Pi}{\partial\theta_j}\right)\right].$$

So, defining $(V_n)_{i,j} = \text{tr}\left(S_n\frac{\partial\Pi}{\partial\theta_i}'H\frac{\partial\Pi}{\partial\theta_j}\right) = \text{tr}\left(\frac{\partial\Pi}{\partial\theta_j}S_n\frac{\partial\Pi}{\partial\theta_i}'H\right)$, we can rewrite the above inequalities as

$$c(\eta)V_n \le B_n \le C(\eta)V_n \tag{39}$$

By the regularity property of the prior distribution we can find a compact $G \subset \Theta_2$ so that $\Sigma(\gamma)$ is nonsingular for $\gamma \in G$ and $P[\gamma \in G] \ge 1 - \eta/2$. Consequently, we can find $h_0, H_0$ so that $h_0I \le H(\gamma) \le H_0I$. Analogous to the proof of (39) above, we can then show that $h_0\Psi_n \le V_n \le H_0\Psi_n$ which, together with (39), proves the lemma. ■

**Lemma A4** $\det(B_n(\beta,\gamma)) = O_P(n^\ell n^p \det(\Psi_n))$.

**Proof** Since $\det(B_n) = \det(B_n^{(\beta)})\det(B_n^{(\gamma)})$, the second proposition of Lemma A1 implies that $\det(B_n^{(\gamma)}) = O(n^p)$. Lemma A3 shows that $\det(B_n^{(\beta)}) = O_P(n^\ell\det(\Psi_n))$, and, the result follows. ■

## 8.2   Proof of Theorem 3

The proof of theorem 3 will take up the remainder of this subsection and will be developed using a series of propositions, whose proofs will be given as we go along and at the end of the subsection. As in theorem 1 and 2, it is helpful to 'cut out' events with small probabilities and in doing so it is convenient to use the notation introduced in Definition 2.

We proceed in an analogous way to the proof of theorem 2. For every $\eta > 0$ we will construct events $C_n = C_n(\eta) \in \mathfrak{F}_n^*$ with $\liminf P(C_n) \ge 1 - \eta$, define the approximating measures $Q_n = Q_n^{(\eta)}$ by $Q_n(A) = P(A \cap C_n)$ and then make use of Lemma RRN to give the density

$$\frac{dQ_n^{(\eta)}}{dP_{(\beta,\gamma)}^{(n)}} = \int_\Theta \mathbf{1}_{C_n}((\kappa,\rho),\cdot)\frac{dP_{(\kappa,\rho)}}{dP_{(\beta,\gamma)}}\pi(\kappa,\rho)d\kappa d\rho.$$

We will show that, on the event $C_n$ (or, to be precise, if $1_{C_n}((\kappa, \rho), \cdot)$ is not identical zero) and using $K_n$ to denote random variables which remain $O_P(1)$,

$$\log \frac{dP_{(\kappa,\rho)}}{dP_{(\beta,\gamma)}} \leq K_n, \tag{40}$$

and, with $\lambda$ denoting the Lebesgue-measure on the appropriate spaces,

$$\lambda(\{\kappa : \mathbf{1}_{C_n}((\kappa, \rho), \cdot) \neq 0\}) \leq \frac{K_n}{\sqrt{n^\ell \det(\Psi_n)}}, \tag{41}$$

$$\lambda(\{\rho : I_{C_n}((\kappa, \rho), \cdot) \neq 0\}) \leq \frac{K_n}{\sqrt{n^p}}. \tag{42}$$

The required result then follows from these bounds.

To start, we write the log likelihood function as

$$\ell(\kappa, \rho) = \tfrac{n}{2} \log \det H(\rho) - \tfrac{1}{2} \sum (y_t - \Pi(\kappa) z_t)' H(\rho)(y_t - \Pi(\kappa) z_t)$$

Setting $u_t = y_t - \Pi(\beta) z_t$, $H_0 = H(\gamma)$, $\Delta(\rho) = \Pi(\beta) - \Pi(\rho)$ we have

$$\begin{aligned}
\ell(\kappa, \rho) - \ell(\beta, \gamma) = {} & \tfrac{n}{2}(\log \det H(\rho) - \log \det H(\gamma)) - \tfrac{1}{2}\mathrm{tr}((\Sigma u_t u_t')(H(\rho) - H(\gamma))) \\
& - \tfrac{1}{2}\Sigma(u_t' H(\rho)\Delta(\rho) z_t + z_t'\Delta(\rho)' H(\rho) u_t) \\
& - \tfrac{1}{2}\Sigma z_t'\Delta(\rho)' H(\rho)\Delta(\rho) z_t
\end{aligned}$$

For (40) to hold, we need to show that this difference in the likelihoods remains bounded in probability. As we only need to give upper bounds for these terms, we only have to deal with the first two summands on the right side. This is accomplished in the two propositions that follow.

### 8.2.1 Proposition A5.

*For every $\eta > 0$ there exists a sequence $C_n^{(1)} = C_n^{(1)}(\eta)$ of events so that $\liminf P(C_n^{(1)}) \geq 1 - \eta$ and the following property holds: if $\mathbf{1}_{C_n^{(1)}}((\kappa, \rho), \cdot)$ is not identical zero, then*

$$\tfrac{n}{2}(\log \det H(\rho) - \log \det H(\gamma)) - \tfrac{1}{2}\mathrm{tr}((\Sigma u_t u_t')(H(\rho) - H(\gamma)) \leq L_n$$

*where the $L_n$ (for each fixed $\eta$) are $O_P(1)$ random variables.*

### 8.2.2 Proposition A6.

*For every $\eta > 0$ there exists a sequence $C_n^{(2)} = C_n^{(2)}(\eta)$ of events so that $\liminf P(C_n^{(2)}) \geq 1 - \eta$ and the following property: if $\mathbf{1}_{C_n^{(2)}}((\kappa, \rho), \cdot)$ is not identical zero, then*

$$\left| -\tfrac{1}{2}\Sigma(u_t' H(\rho)\Delta(\kappa) + \Delta(\kappa)' H(\rho) u_t) \right| \leq L_n$$

26

*where $L_n$ are $O_P(1)$ random variables.*

Analogous to the proof of Theorem 2, we will "cut out" all parameters "far away" from $(\beta, \gamma)$. Consider the OLS-estimator for $\Pi$ and $\Sigma$, viz., $\hat{\Pi}_n = (\Sigma y_t z_t')(\Sigma z_t z_t')^{-1}$ and $\hat{\Sigma}_n = \frac{1}{n}\Sigma(y_t - \hat{\Pi}_n z_t)(y_t - \hat{\Pi}_n z_t)'$.

### 8.2.3 Proposition A7.

$$\sqrt{n}(\hat{\Pi}_n - \Pi(\beta))O_n' D_n,$$

$$\sqrt{n}(\hat{\Sigma}_n - \Sigma(\gamma))$$

*and*

$$\sqrt{n}\left(\hat{\Sigma}_n - \frac{1}{n}\sum u_t u_t'\right)$$

*remain $O_P(1)$ as $n \to \infty$.*

### 8.2.4 Proof of Proposition A7.

Since $\hat{\Pi}_n - \Pi(\beta) = (\Sigma u_t z_t')(\Sigma z_t z_t')^{-1}$, it is easily seen from Assumption **C3** that $\sqrt{n}(\hat{\Pi}_n - \Pi(\beta))O_n' D_n$ converges in distribution to $WA^{-1}$, which proves the first statement in view of Lemma P-BD. For the second, observe that

$$\begin{aligned}
\hat{\Sigma}_n - \Sigma =\ & \left(\tfrac{1}{n}\Sigma u_t u_t' - \Sigma\right) \\
& + 2 \cdot \tfrac{1}{\sqrt{n}}(\hat{\Pi}_n - \Pi(\beta))O_n D_n(D_n^{-1}O_n')\tfrac{1}{\sqrt{n}}\Sigma z_t u_t' \\
& + (\hat{\Pi}_n - \Pi(\beta))O_n D_n\left\{(D_n^{-1}O_n')\left(\tfrac{1}{n}\Sigma z_t z_t'\right)(O_n D_n^{-1})\right\}D_n O_n'(\hat{\Pi}_n - \Pi(\beta))'.
\end{aligned}$$

**C3** and the first result of this lemma now show that the second and the third statements of the lemma hold, again in view of Lemma P_BD. ∎

### 8.2.5 Proof of Proposition A6.

Fix $\eta > 0$. Then we can find $M = M(\eta)$ so that $P(C_n') > 1 - \eta/2$ and $P(C_n'') > 1 - \eta/2$ with $C_n' = [||\sqrt{n}(\hat{\Pi}_n - \Pi)O_n D_n|| < M]$ and $C_n'' = [||H|| < M]$. Define $C_n^{(2)}$ as $C_n' \cap C_n''$. Then

$$\begin{aligned}
& \Sigma(u_t' H(\rho)\Delta(\kappa)z_t + z_t'\Delta(\kappa)' H(\rho)u_t) \\
=\ & 2\,\mathrm{tr}\left(\{D_n^{-1}O_n'(\Sigma z_t u_t')\}H(\rho)\{(((\hat{\Pi}_n - \Pi(\beta))O_n D_n\}\right) \\
& + 2\,\mathrm{tr}\left(\{D_n^{-1}O_n'(\Sigma z_t u_t')\}H(\rho)\{(\Pi(\kappa) - \hat{\Pi}_n)O_n D_n\}\right).
\end{aligned}$$

Now analyse the two summands on the right-hand side of this equation. Each of these is a trace of a product of three (random) matrices. The first factor is a random matrix

which converges in distribution. The norm of the second is, provided $\mathbf{1}_{C_n^{(2)}}((\kappa, \rho), \cdot)$ is not identical zero, dominated by $M$. Due to the construction of $C_n^{(2)}$, the same applies to the third factor of the second sum. The third factor in the first sum is the product of random matrices which converge in distribution to $WA^{-1}$. Applying Lemma P_BD completes the proof. ∎

### 8.2.6  Proof of Proposition A5.

Proposition A5 can be proven in a similar manner. Let $\hat{H}_n = \widehat{\Sigma}_n^{-1}$ and write

$$
\begin{aligned}
&\sqrt{n}\tfrac{1}{2}(\log \det H(\rho) - \log \det H(\gamma)) - \tfrac{1}{2}\operatorname{tr}(\tfrac{1}{n}(\Sigma u_t u_t')(H(\rho) - H(\gamma)) \\
&= \sqrt{n}\left\{\tfrac{1}{2}(\log \det H(\rho) - \log \det \hat{H}_n) - \tfrac{1}{2}\operatorname{tr}(\hat{H}_n^{-1}(H(\rho) - \hat{H}_n))\right\} \\
&\quad + \sqrt{n}\left\{\tfrac{1}{2}(\log \det \hat{H}_n - \log \det H(\gamma)) - \tfrac{1}{2}\operatorname{tr}(\hat{H}_n^{-1}(\hat{H}_n - H(\gamma)))\right\} \\
&\quad + \sqrt{n}\left\{\tfrac{1}{2}\operatorname{tr}((\hat{H}_n^{-1} - \tfrac{1}{n}(\Sigma u_t u_t'))(H(\rho) - \hat{H}_n))\right\} \\
&\quad + \sqrt{n}\left\{\tfrac{1}{2}\operatorname{tr}((\hat{H}_n^{-1} - \tfrac{1}{n}(\Sigma u_t u_t'))(H(\gamma) - \hat{H}_n))\right\}.
\end{aligned}
\tag{43}
$$

Deal with each of the four terms (in braces) on the right-hand side separately. Choose an arbitrary $\eta > 0$. Then we can find $M = M(\eta)$ so that, with $C_n^{(1)} = [||\hat{H}_n - H|| \le M/\sqrt{n}]$, $P(C_n^{(1)}) \ge 1 - \eta$. Then Proposition A7 immediately shows that the fourth term converges to zero and the third term is dominated on $C_n^{(1)}$ by

$$
\sup_{\{\rho:\ \mathbf{1}_{C_n^{(1)}}(\rho, \cdot)) \ne 0\}} \sqrt{n}\left\{\tfrac{1}{2}\operatorname{tr}((\hat{H}_n - \tfrac{1}{n}(\Sigma u_t u_t'))(H(\rho) - \hat{H}_n))\right\} \to 0,
$$

as the first factor within the trace converges to zero from Proposition A7 and the second factor remains bounded.

For the first and second terms of (43) we use the expansion for $\log \det A$ that is given in Proposition A8, stated at the end of this section. This proposition shows that the difference of the second term of (43) and

$$
\operatorname{tr}\left(\sqrt{n}(H(\gamma) - \tfrac{1}{n}(\Sigma u_t u_t'))\hat{H}_n^{-1}\sqrt{n}\left((H(\gamma) - \tfrac{1}{n}(\Sigma u_t u_t'))\right)\hat{H}_n^{-1}\right)
$$

converges in probability to zero. As this sequence obviously converges in distribution, we can apply Lemma P-BD and it remains $O_P(1)$.

Now we only have to analyse the first summand in (43). Using the defining property of $C_n^{(1)}$, it is easily seen that

$$
\sup_{\{\rho:\ \mathbf{1}_{C_n^{(1)}}(\rho, \cdot) \ne 0\}} |h_{1,n}(\rho) - h_{2,n}(\rho)| \to 0,
$$

and

$$
h_{1,n}(\gamma) - h_{2,n}(\gamma) \to 0,
$$

where

$$h_{1,n}(\rho) = \sqrt{n} \left\{ \tfrac{1}{2}(\log \det H(\rho) - \log \det \hat{H}_n) - \tfrac{1}{2}\mathrm{tr}(\hat{H}_n^{-1}(H(\rho) - \hat{H}_n)) \right\},$$

and

$$h_{2,n}(\rho) = \mathrm{tr}(\sqrt{n}(H(\rho) - \hat{H}_n)\hat{H}_n^{-1}\sqrt{n}(H(\rho) - \hat{H}_n)\hat{H}_n^{-1}).$$

It is clear that $h_{2,n}(\gamma)$ converges in distribution and again remains $O_P(1)$ by virtue of Lemma P-BD, so we only have to analyse $h_{2,n}(\rho)$. For doing this, observe that Proposition A7 implies that $\sqrt{n}||\hat{H}_n - H(\gamma)||$ remains $O_P(1)$, and so

$$\sup_{\{\rho:\, \mathbf{I}_{C_n^{(1)}}(\rho,\cdot)\neq 0\}} \sqrt{n}||H(\gamma) - H(\rho)||$$

remains $O_P(1)$, too. We can therefore conclude (with the help of Assumption **C2** on local identification) that

$$s_n = \sup_{\{\rho:\, \mathbf{I}_{C_n^{(1)}}(\rho,\cdot)\neq 0\}} \sqrt{n}||\gamma - \rho|| \tag{44}$$

is $O_P(1)$. Moreover, $||\sqrt{n}(H(\rho) - \hat{H}_n)|| \leq ||\sqrt{n}(H(\gamma) - \hat{H}_n)|| + ||\sqrt{n}(H(\rho) - H(\gamma))||$. The first of these summands remains $O_P(1)$. The second one, if $\mathbf{I}_{C_n^{(1)}}(\rho,\cdot) \neq 0$, is dominated by

$$s_n \sup_{\{\rho:\mathbf{1}_{C_n^{(1)}}(\rho,\cdot)\neq 0\}} ||DH(\rho)||, \tag{45}$$

where $DH = \left( \frac{\partial H}{\partial \gamma_1}, ..., \frac{\partial H}{\partial \gamma_\ell} \right)$ is the matrix composed of the first derivatives. Since for an arbitrary small $\kappa > 0$ $\{\rho : \mathbf{1}_{C_n^{(1)}}(\rho,\cdot) \neq 0\} \subset \{\rho : ||\rho - \gamma|| < \kappa\}$ for all but a finite number of $n$, we may conclude that both factors of our product (45) remain $O_P(1)$. This completes the proof of Proposition A5. ∎

Now, since $\lambda(\{\rho : \sqrt{n}||\gamma - \rho|| \leq s_n\}) = \mathrm{const}\cdot(\sqrt{n})^{-p}s_n^p$, and $s_n$ is $O_P(1)$ from (44) above, we have proved (42).

To complete the mainline of our proof, it remains to show (41). Let us define our events for some given $\eta > 0$. In particular, using Proposition A7 we can find an $M = M(\eta)$ so that $P(C_n^{(2)}) > 1 - \eta$ with $C_n^{(2)} = [||(\sqrt{n}\left(\hat{\Pi}_n - \Pi\right)O_n'D_n)|| \leq M]$. Then, we have to show that

$$\lambda(\{\kappa : \mathbf{I}_{C_n^{(2)}}(\kappa,\cdot) \neq 0\}) = O_P(n^{-\ell/2}/\sqrt{\det B_n}).$$

As

$$||(\sqrt{n}(\hat{\Pi}_n - \Pi(\beta))O_n'D_n)|| + ||(\sqrt{n}(\hat{\Pi}_n - \Pi(\kappa))O_n'D_n)|| \geq ||(\sqrt{n}(\Pi(\kappa) - \Pi(\beta))O_n'D_n)||$$

we may conclude that

$$\{\kappa : \mathbf{I}_{C_n^{(2)}}(\kappa,\cdot) \neq 0\} \subset \{\kappa : ||(\sqrt{n}(\Pi(\kappa) - \Pi(\beta))O_n'D_n)|| \leq 2M\} \tag{46}$$

29

on the event

$$[||(\sqrt{n}(\hat{\Pi}_n - \Pi(\beta))O_n' D_n)|| \leq M]. \tag{47}$$

This should be understood as follows. For all $\omega$ satisfying event (47) $\{\kappa : \mathbf{I}_{C_n^{(2)}}(\kappa, \omega) \neq 0\}$ is a subset of the set on the right-hand side of (46). By the definition of $M$, the probability of the event (47) is greater than $1 - \eta$. We have to show that for the sets

$$R_n(M) = \{\kappa : ||(\sqrt{n}(\Pi(\kappa) - \Pi(\beta))O_n' D_n)|| \leq 2M\},$$

$\lambda(R_n(M))$ has the correct order of magnitude. We will give the proof only for the case of

$$O_n = I \tag{48}$$

Since the $O_n$ have been assumed to be orthogonal, the proof is easily extended to the general case, but the more complicated notation required would distract from the basic intuition behind the proof. Moreover, we will use the notation Const as a generic symbol for a *strictly positive constant* which is not necessarily the same in every expression. This property is most important for the proof. For reasons of brevity, we will refrain from mentioning the *strict positiveness* of Const every time we use the symbol.

Applying Proposition A7, it is sufficient to show, under our simplifying assumption (48, that $\lambda(R_n(M)) = O_P(n^{-\ell/2}/\sqrt{\det B_n})$ for all $M$. Since all norms on finite-dimensional spaces are equivalent, it is easily seen that

$$R_n(M) \subset \{\kappa : \text{tr}((\sqrt{n}(\Pi(\kappa) - \Pi(\beta))D_n)(\sqrt{n}(\Pi(\kappa) - \Pi(\beta)D_n)' \leq \text{const } M^2\}.$$

Moreover, it is an immediate consequence of Lemma A3 that the volume of the ellipsoid $\{\kappa : n(\kappa - \beta)'\Psi_n(\kappa - \beta) \leq \text{const}\}$ is $O_P(n^{-\ell/2}/\sqrt{\det B_n})$.

Therefore, it is sufficient to show that for each $\beta$ there exist a neighbourhood $U(\beta)$ and a constant $\text{Const} = \text{Const}(\beta)$ so that for $\kappa \in U(\beta)$

$$\text{tr}\left((((\Pi(\kappa) - \Pi(\beta))D_n)((\Pi(\kappa) - \Pi(\beta))D_n)')\right) \geq \text{Const} \cdot (\kappa - \beta)'(\Psi_n/n)(\kappa - \beta). \tag{49}$$

Let $\Pi = \left(\pi^{(1)}, ..., \pi^{(\ell)}\right)$ and $D_n = \text{diag}(\lambda_{1,n}, ..., \lambda_{\ell,n})$. Then, the left side of (49) equals

$$\sum \lambda_{j,n}^2 \left\|\pi^{(j)}(\kappa) - \pi^{(j)}(\beta)\right\|^2,$$

and the right hand side is

$$\sum (\kappa - \beta)_i (\kappa - \beta)_j \text{tr}\left(D_n D_n' \frac{\partial \Pi'}{\partial \theta_i} \frac{\partial \Pi}{\partial \theta_j}\right) = \sum \lambda_{j,n}^2 \left\|\sum (\kappa - \beta)_i \frac{\partial \pi^{(j)}}{\partial \beta_i}\right\|^2,$$

where $||v|| = \sqrt{\sum v_i^2}$ is the usual Euclidean norm. Therefore, we prove the proposition if we can show that for all $j$ and all $\beta$ there exists a neighbourhood $U(\beta)$ so that

$$||\pi^{(j)}(\kappa) - \pi^{(j)}(\beta)||^2 \geq \text{Const} \cdot \left\|\sum (\kappa - \beta)_i \frac{\partial \pi^{(j)}}{\partial \beta_i}\right\|^2. \tag{50}$$

30

At the first sight, the proof of this inequality seems to be a standard exercise in elementary analysis, but this is true only in the case where the right-hand side is nonzero for all nontrivial vectors $(\kappa - \beta)$. One does, however, encounter the problem that, in general, there will exist vectors $(\kappa - \beta)$ that annihilate the right hand side (i.e., $\pi^{(j)}(\cdot)$ has a zero derivative in that direction), so the inequality is trivial for them. But what happens "near" these vectors, i.e., when we add a small component of a vector for which the directional derivative is nonzero)? The left hand side of the inequality will be "small" and so the inequality is nontrivial. The key to establishing the inequality in such neighbourhoods lies essentially in "projecting" the mapping to some lower-dimensional manifolds on which it is regular. We make this construction in what follows.

Let us now fix a $j$ and define $R_\pi = \text{span}\left\{\frac{\partial \pi^{(j)}}{\partial \beta_i}\right\}$, i.e. the vector space of all linear combinations of the $\frac{\partial \pi^{(j)}}{\partial \beta_i}$, and let $N = \left\{h : \frac{\partial \pi^{(j)}}{\partial h} = \sum h_i \frac{\partial \pi^{(j)}}{\partial \beta_i} = 0\right\}$. Further, let $V$ be the orthogonal complement of $N$. If $V$ consists only of the null-vector, then the right-hand side of (50) is identically zero and the inequality is trivial. Hence we can assume that $\dim V > 0$. Then it is easily seen that $\dim R_\pi = \dim V = J$. Then we can find vectors $b_1, ..., b_J$ that form a basis of $V$, i.e., they are linearly independent and $V = \{\sum_{i=1}^{J} \nu_i b_i\}$. It can immediately be seen that there exists a linear, bijective mapping $\varphi : \mathbb{R}^J \to R_\pi$ defined by $\varphi((\nu_1, ..., \nu_J)') = \sum_{i=1}^{J} \nu_i b_i$.

Analogously, we can find a basis $c_1, ..., c_J$ of $R_\pi$; Let us now define $P$ as the $J \times \ell$-matrix describing the *orthogonal projection* onto $R_\pi$ with respect to the basis $c_1, ..., c_J$. That is, for any vector $x \in \mathbb{R}^\ell$, the vector $Px \in \mathbb{R}^J$ is such that $\sum (Px)_i c_i$ is the orthogonal projection of $x$ onto $R_\pi$. It is obvious that

$$\text{rank } P = J. \tag{51}$$

Next, let $p(\cdot)$ be the mapping defined on a neighborhood of the origin of $\mathbb{R}^J$ by the following. If $\nu = (\nu_1, ..., \nu_J)$ then

$$p(\nu) = P(\pi^{(j)}(\beta + \sum \nu_i b_i) - \pi^{(j)}(\beta)).$$

In view of Proposition A9, which is stated and proved at the end of this section, this mapping has the property that the Jacobian of $p(\cdot)$ has full rank at the origin so that $\dim R = \dim V$.

Let $S$ be the projection (defined in $\mathbb{R}^\ell$) on $V$ in direction $N$ (i.e., for $h \in N$, $Sh = 0$, for $h \in V$, $Sh = h$). Since $V$ is the orthogonal complement of $N$, $S$ is an orthogonal projection and therefore

$$||h||^2 \geq ||Sh||^2 \tag{52}$$

Furthermore, it is easily seen that there for all $h \in \mathbb{R}^\ell$

$$||Sh||^2 \geq \text{Const} \cdot \left\|\sum h_i \frac{\partial \pi^{(j)}}{\partial \beta_i}\right\|^2, \tag{53}$$

31

where

$$\text{Const} > 0. \tag{54}$$

A feasible choice of Const is $\min_{h \in \Xi} ||Sh||^2$ with $\Xi = \left\{ h \in V : \left\| \sum h_i \frac{\partial \pi^{(j)}}{\partial \beta_i} \right\|^2 = 1 \right\}$.
$\Xi$ is easily seen to be a compact set, so the infimum of a continuous function on the set is its minimum. Hence any strictly positive function can be bounded from below with a constant greater zero, and so this definition of Const fulfils (54)).

As $\pi^{(j)}(\cdot)$ is continuous, there is a neighborhood $U$ around the $\beta$ for which with $\kappa \in U$ we have $P(\pi^{(j)}(\kappa) - \pi^{(j)}(\beta)) \in W$. Let us analyse the mapping $f$ defined by $f(\kappa) = (\psi \circ \varphi^{-1})(P(\pi^{(j)}(\kappa) - \pi^{(j)}(\beta)))$. Then, due to the differentiability of $\psi$ and $\varphi$, there exists a Const with

$$||f(\kappa)|| \le \text{Const} \cdot ||P(\pi^{(j)}(\kappa) - \pi^{(j)}(\beta))||.$$

Again, Const can be assumed to be $> 0$ without limitation in generality, so we have also

$$\text{Const} \cdot ||f(\kappa)|| \le .||P(\pi^{(j)}(\kappa) - \pi^{(j)}(\beta))||.$$

Then we have for $\kappa \in U$

$$||\pi^{(j)}(\kappa) - \pi^{(j)}(\beta)||^2 \ge ||P(\pi^{(j)}(\kappa) - \pi^{(j)}(\beta))||^2 \ge \text{Const}||f(\kappa)||^2$$

Now it remains to show that $||f(\kappa)||^2 \ge \text{Const} \cdot \left\| \sum (\kappa - \beta)_i \frac{\partial \pi^{(j)}}{\partial \beta_i} \right\|^2$. To prove this inequality it is (because of (53) and (52)) sufficient to show that $||Sf(\kappa)||^2 \ge \text{Const}||S(\kappa - \beta)||^2$ for $||\kappa - \beta||$ sufficiently small. Denoting by $Df$ the Jacobian of $f$, $Sf = \int_0^1 SDf(\beta + \lambda(\kappa - \beta)) \cdot (\kappa - \beta)d\lambda$ we have

$$
\begin{aligned}
&||Sf(\kappa)||^2 \\
&= \int_0^1 \int_0^1 (SDf(\beta + \lambda(\kappa - \beta)) \cdot (\kappa - \beta))'(SDf(\beta + \mu(\kappa - \beta)) \cdot (\kappa - \beta))d\lambda d\mu \\
&= \int_0^1 \int_0^1 (S(\kappa - \beta))'(Df(..))'(Df(..))(S(\kappa - \beta))d\lambda d\mu.
\end{aligned}
\tag{55}
$$

By the chain rule, the Jacobian is

$$Df = (D\psi)(D\varphi)^{-1}P\frac{\partial \pi^{(j)}}{\partial \beta}. \tag{56}$$

Therefore, due to the continuity of $Df$, we can, for $||\kappa - \beta||$ sufficiently small, conclude that

$$||(Df(..))'(Df(..)) - (Df(\beta))'(Df(\beta))||$$

can be made arbitrarily small. Therefore, there exists a neighbourhood around $\beta$ so that the difference for all $\kappa$ from this neighbourhood is less than

$$\lambda_0 = \tfrac{1}{2} \min_{\{h \in V : ||h||^2 = 1\}} h'(Df(\beta))'(Df(\beta))h, \tag{57}$$

32

which is nonzero due to Proposition A10, which is stated and proved below. Thus, for $\kappa$ from this neighbourhood, we can conclude that the integrand in (55) is $\geq \frac{1}{2}\lambda_0||S(\kappa-\beta)||^2$, which completes the proof of (41) and concludes the proof of theorem 3.

To complete the reasoning, it remains only to prove the following three propositions that were used in the proof of Theorem 3.

### 8.2.7   Proposition A8.

*Let $A$, $B$ be nonnegative definite matrices so that $||A - B||\,||B^{-1}|| < 1$. Then*

$$\log \det A - \log \det B - \operatorname{tr}(B^{-1}(A - B)) \tag{58}$$
$$= \operatorname{tr}\left((A - B)\,B^{-1}\,(A - B)\,B^{-1}\right) + o\left(\frac{||B^{-1}||^3||A - B||^3}{1 - ||B^{-1}||\,||A - B||}\right).$$

### 8.2.8   Proof of Proposition A8

This is based simply on a Taylor expansion of $\log \det A$ and is omitted.

### 8.2.9   Proposition A9

*The Jacobian of $p(\cdot)$ has full rank at the origin, namely $\dim R = \dim V$.*

### 8.2.10   Proof of Proposition A9

Assume otherwise. Then, we would be able to find nontrivial $\gamma_i$ so that $\sum \gamma_i \frac{\partial p}{\partial \nu_i} = P\left(\sum \gamma_i \frac{\partial \pi^{(j)}}{\partial b_i}\right) = 0$. By definition of $R_\pi$, $\left(\sum \gamma_i \frac{\partial \pi^{(j)}}{\partial b_i}\right) \in R_\pi$, so if the orthogonal projection of this vector is zero, the vector is zero itself. So we may conclude that $\sum \gamma_i \frac{\partial \pi^{(j)}}{\partial b_i} = 0$ and therefore, since we assumed the $\gamma_i$ to be nontrivial,

$$\frac{\partial \pi^{(j)}}{\partial \gamma} = 0 \text{ with } \gamma = \sum \gamma_i b_i \in R_\pi.$$

But this would imply that $\gamma \in N$, so $\gamma \in R_\pi \cap N = \{\mathbf{0}\}$, which contradicts our assumption of $\gamma$ being nontrivial. Then standard analysis shows that there exists an open set $W \subset \mathbb{R}^J$ around the origin for which there exists an inverse function $\psi$. It is easily seen to be continuously differentiable, and its Jacobian has full rank, too, i.e.,

$$\text{Rank } D\psi = J, \tag{59}$$

giving the required result. ∎

### 8.2.11 Proposition A10

*Let $\lambda_0$ be as defined in (57): Then $\lambda_0 > 0$.*

**Proof.** First observe that the set $\{h \in V : ||h||^2 = 1\}$ is compact: Therefore the *infimum* of a continuous function over this set is a *minimum.* Therefore. our definition in (57) makes sense and we can assume that there exists a $h \in V$ with $||h||^2 = 1$ so that $h'(Df(\beta))'(Df(\beta))h = \lambda_0$. Now suppose the proposition does not hold and there exists $h \in V$ with $||h||^2 = 1$ for which $h'(Df(\beta))'(Df(\beta))h = 0$. Then, $(Df(\beta))h = 0$ and, due to (56) and the nonsingularity of $D\psi(\beta)$ and $D\varphi$, we may conclude that $P\frac{\partial \pi^{(j)}}{\partial \beta}h = P\frac{\partial \pi^{(j)}}{\partial h} = 0$. Since $P$ describes the orthogonal projection onto $R_\pi = \mathrm{span}\left\{\frac{\partial \pi^{(j)}}{\partial \beta_i}\right\} \ni \frac{\partial \pi^{(j)}}{\partial h}$, we may conclude that $\frac{\partial \pi^{(j)}}{\partial h} = 0$. But this is just the definition of $h \in N$ and therefore we have a contradiction to our assumptions (viz., that $h$ was nontrivial and an element of $V$, which defined as the orthogonal complement of $N$). ∎

## 8.3  Proof of Theorem 4

Fix an arbitrary predictor $\widehat{y}_t$. Then, the conditional probability measures $G(\widehat{y}_t, \Sigma)$ (for $t \geq 1$, our common point of initialization for the predictions) produce an empirical model in the sense of earlier sections. One can easily see that the corresponding log likelihood ratio with respect to the true model is esentially given by $\Delta_t$. Therefore, it is apparent that Theorem 4 is a simple consequence of Theorem 1 if we can prove that

$$\frac{\log \det B_n}{\log n} \to K, \tag{60}$$

in probability for Lebesgue-almost all parameters.

We start by choosing arbitrary matrices $B$ and $\Gamma$ that fulfill our idenfication requirements given in Assumption **D1**. Next, we proceed to compute the information matrix $B_n$. Lemma A1 shows that this matrix is block diagonal. To use Lemma A1, it helps to simplify some formulae by defining a $2\ell$-vector $x_t^*$ as follows. The first $\ell$ components of $x_t^*$ are set to the vector $x_t$ itself, and components $\ell + 1$ to $2\ell$ are set equal $-\Pi x_t$. The process $x_t^*$ helps to simplify the expression for the score. For this purpose, define for $(i,j) \in M_1$ the elements of a $k \times 2\ell$ selector matrix $P_{i,j}^1$ to be all zero except the element in position $(i,j)$, which is set to unity. Analogously, define for $(i,j) \in M_2$ the elements of the $k \times 2\ell$ matrix $P_{i,j}^2$ to be zero except the element in $(\ell + i, j)$, which is set to unity. In general, we will write $P_{i,j}^a$ with $a = 1, 2$ corresponding to the indices of $M_1$ and $M_2$, respectively.

We need the following expressions for the derivative matrices: first,

$$\frac{\partial \Pi}{\partial \Gamma_{i,j}} = - \left[ \underset{\text{column } j}{0, 0, ..., (\Gamma^{-1})_i, ..., 0} \right] \Pi,$$

where the first matrix in this product is square and has the $i$-th column of $\Gamma^{-1}$ as its $j$-th column, and zeros elsewhere; and, second,

$$\frac{\partial \Pi}{\partial B_{i,j}} = \left[ 0, 0, .., \underset{\text{column } j}{(\Gamma^{-1})_i}, .., 0 \right].$$

Using the selector matrices $P_{i,j}^a$, we can write these matrices in the form

$$\frac{\partial \Pi}{\partial \Gamma_{i,j}} x_t = \Gamma^{-1} P_{i,j}^1 x_t^* \tag{61}$$

and

$$\frac{\partial \Pi}{\partial B_{i,j}} x_t = \Gamma^{-1} P_{i,j}^2 x_t^* \tag{62}$$

We now proceed to compute the matrix $B_n$. We can think of $B_n$ as a matrix indexed with pairs of elements of $M_a$, which constitute triples when combined with the index $a$. Formulae (61) and (62) allow us to apply Lemma A1 and with a bit of calculation it is readily seen that

$$(B_n)_{(i,j,b),(q,\ell.d)} = \sum_{t \le n} \text{tr}(x_t^*(x_t^*)' P_{i,j}^{b'} \Gamma^{-1} \Sigma \Gamma^{-1} P_{q,\ell}^d) = \sum_{t \le n} \text{tr}(P_{q,\ell}^d x_t^*(x_t^*)' P_{i,j}^{b'} \Gamma^{-1} \Sigma \Gamma^{-1}).$$

For each invertible $\Gamma$ we can find $\delta_1 = \delta_1(\Gamma, \Sigma)$, $\delta_2 = \delta_2(\Gamma, \Sigma)$ so that $\delta_2 I \ge \Gamma^{-1} \Sigma \Gamma^{-1} \ge \delta_1 I$, where $I$ is the identity matrix. Define the matrices $R_n$ by

$$(R_n)_{(i,j,b)(q,\ell,d)} = \sum_{t \le n} \text{tr}(P_{q,\ell}^d x_t^*(x_t^*)' P_{i,j}^{b'}) = \sum_{t \le n} \text{tr}(x_t^*(x_t^*)' P_{i,j}^{b'} P_{q,\ell}^d).$$

Lemma A2 implies that

$$\delta_1 R_n \le B_n \le \delta_2 R_n.$$

Hence, (60) is equivalent to

$$\frac{\log \det R_n}{\log n} \to K. \tag{63}$$

Let us now look at the elements $(R_n)_{(i,j,b)(q,\ell,d)}$ if $i \ne q$. In this case it is easily seen that $P_{i,j}^{b'} P_{q,\ell}^d = 0$ and, therefore,

$$(R_n)_{(i,j,b)(q,\ell,d)} = 0 \text{ for } i \ne q \tag{64}$$

Let us for $1 \le i \le k$ define the matrices $(R_n^{(i)})_{(j,b)(\ell,d)} = (R_n)_{(i,j,b)(i,\ell,d)}$, where $j \in M_1(i)$ if $b = 1$ and $j \in M_2(i)$ if $b = 2$. Then (64) shows that by reordering rows and columns we can rearrange the matrix $R_n$ into the form

$$R_n = \begin{pmatrix} R_n^{(1)} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & R_n^{(k)} \end{pmatrix},$$

which implies that $\det R_n = \prod_{i=1}^n \det R_n^{(i)}$ and consequently

$$\log \det R_n = \sum_{i=1}^n \log \det R_n^{(i)}.$$

Consequently, to prove (63) it is sufficient to show that

$$\frac{\log \det R_n^{(i)}}{\log n} \to TI(r_t^{(i)}),$$

where the processes $r_t^{(i)}$ were defined above. To do so, we have to analyze the matrices $R_n^{(i)}$. Fix $i$, with $1 \le i \le k$, and examine the vector $(P_{i,j}^a x_t^*)$, where $j$ is from $M_a(i)$. Then its components are zero except for the $j$-th, which equals $(x_t^*)_j$: Hence, if $j$ and $\ell$ are from $M_a(i)$, then

$$(R_n^{(i)})_{(j,a)(\ell,a)} = \sum_{t \le n} \mathrm{tr}(P_{i,\ell}^a x_t^*(x_t^*)' P_{i,j}^{a\prime}) = \sum_{t \le n} (x_t^*)_j (x_t^*)_\ell.$$

Going back to the definition of $r_t^{(i)}$, it is apparent that the components of this vector coincide with the components of $x_t^*$ when the index is in $M_1(i)$ or the difference of the index with $m$ is in $M_2(i)$. Both vectors simply pick off the components of $x_t^*$ which correspond to unknown parameters in row $i$. Therefore $R_n^{(i)}$ is just a reodered form of $\sum_{t \le n} r_t^{(i)} r_t^{(i)\prime}$: We therefore have to show that

$$\frac{\log \det \sum_{t \le n} r_t^{(i)} r_t^{(i)\prime}}{\log n} \to TI(r_t^{(i)}). \tag{65}$$

To establish (65), it will be sufficient to prove the following two results:
(i) the existence of diagonal matrices $D_{in}$ and a nonsingular matrix $A_i$ so that

$$D_{in}^{-1} A_i \sum_{i \le n} r_t^{(i)} r_t^{(i)\prime} A_i D_{in}^{-1} \Rightarrow C_i, \tag{66}$$

where $\Rightarrow$ denotes weak convergence and $C_i$ is a (possibly random) matrix which is a.s. invertible; and
(ii)

$$\frac{\log \det D_{in}}{\log n} = TI(r_t^{(i)})/2. \tag{67}$$

We will not give an explicit formula for $A_i$ in (66). We will show its existence, mainly by using permutations and linear combinations of the components of $r_t^{(i)}$ that are analogous to the Gaussian elimination algorithm for solving linear equations. Let us assume our vector has $n_{\text{stat}}$ stationary components, and that there are $n_{\text{coint}}$ linear independent cointegrating relationships. Using a permutation matrix to rearrange the stationary components and then multiplying by a matrix which performs the

36

cointegrating space mapping, we can construct a nonsingular matrix $A_1$ with the following properties: the last $n_{\text{coint}} + n_{\text{stat}}$ components of the random vector $\rho_t = A_1 r_t^{(i)}$ are stationary processes and the first $(m - ((n_{\text{coint}} + n_{\text{stat}})))$ are nonstationary and, moreover, every linear combination of them is nonstationary, so they are what we call full rank nonstationary. Next we will deal with deterministic trends. We will assume here that only linear trends are included, extensions to higher order polynomial trends being straightforward. Without limitation in generality, we can arrange for this type of trend to occur in the first component (otherwise, simply multiply by a permutation matrix to accomplish this positioning of the elements). So, let us assume that $(\rho_t)_1 = at + W_1(t)$, where $W_1(t)$ contains no deterministic trend. Now multiply $\rho_t$ with a matrix $A_2$ constructed in the following way: row 1 of $A_2$ should be the first row of the identity-matrix; row $i$ should be the $i$-th row of the identity matrix if $(\rho_t)_i$ does not contain a deterministic trend; otherwise assume that $(\rho_t)_i = bt + W_i(t)$, where $W_i(t)$ does not contain a deterministic trend and then the $i$-th row should consist of $(-b/a)$ in the first column (to eliminate the trend in the $i$'th row), 1 in the $i$-th column and 0 in the remaining columns; for $i > (m - (n_{\text{coint}} + n_{\text{stat}}))$ let the $i$-th row of $A_2$ be identical to the $i$-th row of the identity matrix. Next, let $R_t = A_2 \rho_t$. Since $R_t = (A_2 A_1) r_t^{(i)}$ is a linear combination of $r_t^{(i)}$ and the matrix $A_2 A_1$ is nonsingular, it is sufficient to prove the assertions (66) and (67) for $R_t$. We now do so for a 'generic' equation in the system and to simplify formulae simply drop the affix $i$ in our remaining derivations.

Let us define for the case where one element contains a (linear) deterministic trend the diagonal matrix

$$D_n = \text{diag}(n^{3/2}, n^1, ..., n^1, n^{1/2}, ..., n^{1/2})$$

where $m - 1 - (n_{\text{coint}} + n_{\text{stat}})$ diagonal elements equal $n$ and $(n_{\text{coint}} + n_{\text{stat}})$ elements equal $n^{1/2}$. In the case where none of the processes contains a deterministic trend we define

$$D_n = \text{diag}(n^1, ..., n^1, n^{1/2}, ..., n^{1/2}) \tag{68}$$

where the first $(m - (n_{\text{coint}} + n_{\text{stat}}))$ diagonal elements equal $n$ and the rest equal $n^{1/2}$. Now it is easily seen that (67) holds true for our choice of $D_n$. It now remains to show (66): We have to compute the limiting distribution of $D_n^{-1} \sum_{t \leq n} R_t R_t' D_n^{-1}$: Keeping in mind that the vector $R_t$ is composed of linear combinations of the original vector we can apply the limit theory (37) that follows from Assumption **D3.** We will only deal with the case where a linear deterministic trend is present, because the other case follows in an analogous fashion. So, in this case, the first component of $R_n$ contains a deterministic trend and we can partition the vector $R_n$ into three parts. The first part consists of the first component only, the second part comprises the $m - 1 - (n_{\text{coint}} + n_{\text{stat}})$ nonstationary components and the third part consists of the $n_{\text{coint}} + n_{\text{stat}}$ stationary components. Next, we partition the matrices $D_n^{-1} \sum_{t \leq n} R_t R_t' D_n^{-1}$ and their limit random matrices analogously into nine submatrices, so that we have,

37

in effect, to show that

$$D_n^{-1} \sum_{t \leq n} R_t R_t' D_n^{-1} \Rightarrow C = \begin{pmatrix} c_{11} & c_{12} & c_{13} \\ c_{12}' & C_{22} & C_{23} \\ c_{13}' & C_{23}' & C_{33} \end{pmatrix} \tag{69}$$

and the limit matrix $C$ is nonsingular a.s.

We know from the construction of $R_n$ that its *first* component consists of a deterministic trend (plus terms that are of smaller order than n): We therefore may conclude that for $0 \leq z \leq 1$

$$\lim_{n \to \infty} \frac{(R_{nz})_1}{n} \to az$$

and

$$a \neq 0$$

if $a = 0$ no deterministic trend would be present. Therefore, it is easy to see that (69) holds true for the its uppermost left corner with

$$c_{11} = \frac{a^2}{3} = a^2 \int_0^1 z^2 dz.$$

Since the components of $R_n$ are linear combinations of the $z_t$, we can apply Assumption **D3** and (37). In particular, the vector $R_n^{(2)}$ consisting of the *nonstationary* components (i.e., components $2 : (m - (n_{\text{coint}} + n_{\text{stat}})))$ satisfies an invariance principle. There exists a (vector) nonsingular Wiener process $V$ for which with $0 \leq z \leq 1$

$$R_{nz}^{(2)} \Rightarrow V(z) \quad \text{as } n \to \infty,$$

where the convergence is understood as convergence in the Skorohod topology in the function space $D[0,1]$. For the *stationary* components $R_n^{(3)}$ (the remainder of the vector $R_n$) we postulated (among other things) the existence of second moments and ergodicity. Hence, we may conclude that

$$\frac{1}{n} \sum_{i \leq n} R_i \to_{a.s.} \overline{R}$$

and

$$\frac{1}{n} \sum_{i \leq n} R_i R_i' \to_{a.s.} \overline{C},$$

where $\overline{C}$ is nonsingular and $\overline{C} - \overline{R}\overline{R}'$ is nonnegative definite.

Some lengthy calculations which are similar to those in Park and Phillips (1988,1989) and which we therefore omit here, show that (69) is indeed true and we have the following limits

$$c_{12} = \int_0^1 V'(z) z dz,$$

$$c_{13} = \frac{1}{2} a \overline{R}' = \int_0^1 \overline{R}' z dz,$$

38

$$C_{22} = \int_0^1 V(z)V'(z)dz,$$
$$C_{23} = \int_0^1 \overline{R}'V(z)dz,$$
$$C_{33} = \overline{C}.$$

Therefore, it remains to show the nonsingularity of the matrix $C$. Assume the opposite to be true. Then, there exists a vector $d$ for which

$$d'Cd = 0$$

and, using the above expressions, there would exist constants $A$ and vectors $D, E$, not all zero, for which

$$\int_0^1 (Az + D'V(z) + E'R)^2 dz + E'((\overline{C} - \overline{RR}')E = 0$$

Keeping in mind that $\overline{C} - \overline{RR}'$ is nonnegative definite, this would imply that

$$\int_0^1 (Az + D'V(z) + E'R)^2 dz = 0,$$

which obviously contradicts the nonsingularity of the process $V$, so the singularity of C must be wrong.

The proof is now completed for the case where the process contains a linear deterministic trend. If such a trend is not present in the predetermined variables and we have to use (68) in the definition of $D_n$, we can proceed in an analogous manner. The arguments carry over almost verbatim, and one only has to ignore all statements regarding the first component. In a similar way, when there are deterministic trends of degree $p$ in the process, the same arguments apply with the modified definition

$$D_n = \mathrm{diag}(n^{p+\frac{1}{2}}, n^1, ..., n^1, n^{\frac{1}{2}}, ..., n^{\frac{1}{2}})$$

of the normalizing matrix.

## 9   Notation

| | |
|---|---|
| $\to_{a.s.}$ | almost sure convergence |
| $\to_{P_\theta}$ | convergence in $P_\theta$ probability |
| $\Rightarrow, \to_d$ | weak convergence |
| $o_{P_\theta}(1)$ | tends to zero in $P_\theta$ probability |
| $O_{P_\theta}(1)$ | bounded in $P_\theta$ probability |
| $O_P(1)$ | bounded in $P$ probability |
| $\sim_d$ | asymptotically distributed as |
| $I_A(\cdot)$ | indicator function of $A$ |
| $E_\theta$ | expectation under $P_\theta$ |
| $TV(P,Q) = \sup_{A \in \mathfrak{F}} |P(A) - Q(A)|$ | total variation |
| $\lambda_{\min}(B)$ | smallest eigenvalue of $B$ |
| $\|\cdot\|$ | Euclidean norm in $\mathbf{R}^k$ |

# 10 References

## References

[1] Dawid, A. P. (1984). "Present position and potential developments: Some personal views, statistical theory, the prequential approach," *Journal of the Royal Statistical Society*, Series A, 147, 278–292.

[2] Doan, T., R.B. Litterman and C. Sims (1984). "Forecasting and conditional projections using realistic prior distributions," *Econometrics Reviews* 3, 1–100.

[3] Engle, R. F., D. F. Hendry and J. F. Richard (1983): "Exogenity," *Econometrica,* 51, 277–304.

[4] Gerencser, L., J. Rissanen (1992): "Asymptotics of Predictive Stochastic Complexity." In Brillinger, Caines, Geweke, Parzen, Rosenblatt and Taqqu (eds.), *New Directions in Time Series 2.* Springer Verlag, New York, pp. 93–112.

[5] Keuzenkamp, H. A., M. McAleer and A. Zellner (1999). "*Simplicity, Inference and Econometric Modelling*"; Cambridge: Cambridge University Press.

[6] Kim, J. Y. (1994): "Bayesian Asymptotic Theory in a Time Series Model with a Possible Nonstationary Process," *Econometric Theory,* 10, 764–773.

[7] LeCam, L. (1986). *Asymptotic Methods in Statistical Decision Theory.* New York: Springer.

[8] Park, J. Y. and P. C. B. Phillips (1988): "Statistical Inference in Regressions with Integrated Processes: Part 1," *Econometric Theory,* 4, 468–497.

[9] Park, J. Y. and P. C. B. Phillips (1989): "Statistical Inference in Regressions with Integrated Processes: Part 2," *Econometric Theory,* 5, 95–131.

[10] Phillips, P. C. B (1996): "Econometric Model Determination," *Econometrica,* , 64, 763-812.

[11] Phillips, P. C. B. and S. N. Durlauf (1986). "Multiple time series regression with integrated processes," *Review of Economic Studies* 53, 473–496.

[12] Phillips, P. C. B and Werner Ploberger (1992): "Time Series Modeling with a Bayesian Frame of Reference: Concepts, Illustrations and Asymptotics," Cowles Foundation Discussion paper No. 980.

[13] Phillips, P. C. B. and Werner Ploberger (1996): "An Asymptotic Theory of Bayesian Inference for Time Series," *Econometrica,* 64, 381-413.

[14] Phillips, P. C. B. and V. Solo (1992). "Asymptotics for linear processes," *Annals of Statistics* 20, 971–1001.

[15] Rissanen, J. J. (1986): "Stochastic Complexity and Modelling," *Annals of Statistic,* 14, 1080–1100.

[16] Rissanen, J. J. (1987): "Stochastic Complexity" (with discussion), *Journal of the Royal Statistical Society,* 49, 223–239, and 252–265.

[17] Rissanen, J. J. (1996): "Fisher Information and Stochastic Complexity", *IEEE Transactions on Information Theory,* 42, 40–47.

[18] West, M. and P. J. Harrison (1989). *Bayesian Forecasting and Dynamic Models.* New York: Springer–Verlag.

[19] Zellner, A. and C–K. Min (1992). "Bayesian analysis, model selection and prediction," University of Chicago, mimeographed.