# Focused information criterion for locally misspecified vector autoregressive models

Jan Lohmeyer, Franz Palm, Hanno Reuvers & Jean-Pierre Urbain

Taylor & Francis
Taylor & Francis Group

# Focused information criterion for locally misspecified vector autoregressive models

Jan Lohmeyer, Franz Palm, Hanno Reuvers, and Jean-Pierre Urbain*

Department of Quantitative Economics, Maastricht University SBE, Maastricht, The Netherlands

**ABSTRACT**

This paper investigates the focused information criterion and plug-in average for vector autoregressive models with local-to-zero misspecification. These methods have the advantage of focusing on a quantity of interest rather than aiming at overall model fit. Any (sufficiently regular) function of the parameters can be used as a quantity of interest. We determine the asymptotic properties and elaborate on the role of the locally misspecified parameters. In particular, we show that the inability to consistently estimate locally misspecified parameters translates into suboptimal selection and averaging. We apply this framework to impulse response analysis. A Monte Carlo simulation study supports our claims.

## 1. Introduction

The motivation for this paper stems from Hansen (2005). The author considers a Gausssian ARMA(1,1) model approximated by AR($k$) models with $k \in \{0, 1, \ldots, k_{max}\}$ and is interested in the impulse responses. Table 1 of Hansen (2005) shows that the MSE-minimizing AR order depends strongly on parameter values and impulse response horizon. An extreme case is the specification $y_t = 0.5y_{t-1} + \epsilon_t - 0.9\epsilon_{t-1}$. The MSE-minimizing autoregressive orders equal 0 and 10 for the impulse responses at horizon 2 and 6, respectively. Ivanov and Kilian (2005) report a similar issue in a VAR setting. They simulate VAR processes similar to those often found in empirical work, and rank different model selection criteria (AIC, BIC, HQ, and serial correlation tests) based on the MSE of the estimated impulse responses implied by the selected model. A uniformly best criterion was not found. This might be expected since information criteria like AIC and BIC aim at global model fit and do not take into account the quantity of interest (such as the impulse response at a particular horizon). The Focused Information Criteria introduced by Claeskens and Hjort (2003) does take into account the interest of the researcher. Hansen (2005) acknowledged the opportunities for the FIC for the estimation of impulse responses when he remarked based on simulation outcomes: "The message from Tables 2 and 3 is that the FIC is an intriguing challenger to existing model selection methods and deserves attention and scrutiny." A theoretical justification of these simulation results was not provided.

We develop a theoretical framework starting from a vector autoregression where part of the coefficients are local-to-zero, i.e., declining to zero at a rate of $T^{-1/2}$ with $T$ denoting sample size. This case

is fundamentally different from a static setup because dynamic properties are also varying with sample size.[1] Building on ideas from Claeskens and Hjort (2003), Claeskens and Hjort (2008) and Liu (2015), we propose an estimator that can be used for both model selection and model averaging. This estimator is fairly general as it only requires the parameters of interest to be a sufficiently smooth transformation of the model's parameters. The results are subsequently applied to the specific case of impulse responses. A slight generalization of a theorem by Liu (2015) enables us to not only construct confidence intervals for a specific horizon, but to also construct confidence bands for multiple horizons. In addition, we provide an in depth discussion on the role of the local-to-zero parameters. These parameters cannot be estimated consistently, and we show that as a consequence the FIC and plug-in averages do not fully minimize the asymptotic mean squared error.

Our paper is related to the literature on model selection/averaging and the literature on impulse response analysis. We now discuss both. One of the earliest references on frequentist model averaging is the paper by Bates and Granger (1969). The literature on model averaging that is unrelated to forecasting is of a more recent origin. One literature branch on frequentist model averaging started with the paper by Hansen (2007). This paper shows that weight selection by minimization of Mallow's criterion will asymptotically lead to the lowest possible squared error among a class of estimators. His regression setup with homoskedastic errors was generalized to regression forecasts in Hansen (2008), and was modified by Hansen and Racine (2012) to allow for heteroskedastic errors. A time series application to stationary autoregressions of infinite order is Zhang et al. (2013). Zhang and Liu (2017) report results on the distribution of Mallow's and Jackknife-based model averaging weights in linear regressions with irrelevant variables.

The second branch of literature on frequentist model averaging evolves around locally misspecified models. The FIC was proposed by Claeskens and Hjort (2003) and extended by Claeskens and Hjort (2008). The underlying idea has been applied to various settings. We will report a nonexhaustive list. Liu (2015) considered the linear regression setup and derived asymptotically valid confidence intervals. Two additions to the treatment effects literature are Lu (2015) and Kitagawa and Muris (2016). DiTraglia (2016) provides results for generalized method of moments estimation. Liu and Kuo (2016) consider predictive regressions.

Finally, we briefly discuss the literature on impulse response (IR) analysis in autoregressive models. A comprehensive discussion on IR analysis can be found in Section 3.7 of Lütkepohl (2005). Both Lütkepohl (1990) and Benkwitz et al. (2000) have reported that the coverage of the impulse response confidence intervals can be low since the convergence rate of the estimators to their asymptotic distribution is nonconstant over the whole parameter space. Another important topic for impulse response analysis is the construction of joint confidence bands. A naive Cartesian product of the individual confidence intervals leads to severe undercoverage, whereas confidence bands based on the Bonferroni inequality have good coverage but are at the same time excessively wide. Often considered alternatives are bootstrap methods, e.g., Kilian (2001), Lütkepohl et al. (2015), and Bruder and Wolf (2017).

The remainder of this paper is organized as follows. Section 2 presents the model framework, the estimation procedure, and the asymptotic properties of: (1) the parameter estimates, (2) the feasible FIC, and (3) the elements of the weighting matrix. A discussion and illustration of the consequences of the inconsistent estimation of the local-to-zero parameter follows. Our theoretical findings are subsequently supported by various Monte Carlo simulations in Section 3. Section 4 concludes, and the mathematical proofs are presented in the Appendix.

In terms of notation, we follow Abadir and Magnus (2002) as closely as possible; in particular $\xrightarrow{d}$ and $\xrightarrow{p}$ signify convergence in distribution and convergence in probability, respectively. The stochastic and the strict stochastic order relations are denoted by $O_p(\cdot)$ and $o_p(\cdot)$. Vectors are printed in bold and

---

[1]Dynamic models under local-to-zero misspecification were discussed in Claeskens et al. (2007) and Rohan and Ramanathan (2011). Both papers first derive the asymptotic results in a setting without local misspecification and subsequently introduce the misspecification (see p. 363 of Claeskens et al. (2007) and Equation (8) on p. 221 of Rohan and Ramanathan (2011)). The theoretical implications of this two step procedure are not completely clear.

denote column vectors by default. $\mathbf{0}_j$ is the column vector of length $j$ consisting of zeros only. We permit small deviations from Abadir and Magnus (2002) to keep our notation in line with the notation of Liu (2015).

## 2. Theory

### 2.1. Model framework

Let the $K$-dimensional multiple time series $\{\{y_{T,t}\}_{t=-\infty}^{\infty}\}_{T=1}^{\infty}$ constitute a vector triangular array generated by the vector autoregressive (VAR) processes

$$y_{T,t} = B_1 y_{T,t-1} + \ldots + B_{p_1} y_{T,t-p_1} + \frac{\Delta_1}{\sqrt{T}} y_{T,t-p_1-1} + \ldots + \frac{\Delta_{p_2}}{\sqrt{T}} y_{T,t-p_1-p_2} + u_t, \quad (2.1)$$

where $B_i$ ($i \in \{1,2,\ldots,p_1\}$) and $\Delta_i$ ($i \in \{1,2,\ldots,p_2\}$) are ($K \times K$) coefficient matrices. Equation (2.1) differs in one important aspect from the usual VAR specifications, namely some of the coefficient matrices are premultiplied by $T^{-1/2}$ with $T$ denoting sample size. This local-to-zero misspecification causes different dynamics for every $T$. Mathematically, this decay rate will prove to be crucial for the development of the asymptotic theory because it prevents the omitted variable bias from diverging with increasing sample size. Intuitively, we could think of this model specification as expressing a degree of uncertainty concerning the true lag order. The VAR process includes $p := p_1 + p_2$ lags for finite $T$, yet asymptotically a VAR($p_1$) remains. This can be interpreted as exploring a shrinking neighborhood of the VAR($p_1$) model.

### 2.2. Parameter estimation and asymptotics

To simplify notation, we collect all the parameters in the matrices $B = [B_1 \ B_2 \ \cdots \ B_{p_1}]$, $C_T = [\Delta_1 \ \Delta_2 \ \cdots \ \Delta_{p_2}]/\sqrt{T} = \Delta/\sqrt{T}$, and define $\Theta_T = [B \ C_T]$. Similarly to Lütkepohl (2005), we also stack the observations over time to obtain,

$$
\begin{aligned}
Y_T &:= (y_{T,1}, y_{T,2}, \ldots, y_{T,T}) & (K \times T), \\
z_{T,t} &:= \begin{bmatrix} y_{T,t} \\ y_{T,t-1} \\ \vdots \\ y_{T,t-p+1} \end{bmatrix} & (Kp \times 1), \\
Z_T &:= (z_{T,0}, z_{T,1}, \ldots, z_{T,T-1}) & (Kp \times T), \\
U &:= (u_1, u_2, \ldots, u_T) & (K \times T).
\end{aligned}
\quad (2.2)
$$

The model can now be expressed as $Y_T = \Theta_T Z_T + U$. A variety of approximating models can be considered but we will restrict our attention to models that use the same lag order in every equation (see Remark 1 for further details). Using the same lag order in the cross-section is common practice and will decrease the notational burden. Selection matrices are used to relate all estimators to the estimator using $p$ lags. That is, for some integer $m$, such that $p_1 \leq m \leq p$,

$$
\begin{aligned}
L &:= L^{(1)} \otimes I_K, \text{ with } L^{(1)} = \begin{bmatrix} I_{p_1} \\ O_{p_2 \times p_1} \end{bmatrix} & (Kp \times Kp_1), \\
S_0 &:= S_0^{(1)} \otimes I_K, \text{ with } S_0^{(1)} = \begin{bmatrix} O_{p_1 \times p_2} \\ I_{p_2} \end{bmatrix} & (Kp \times Kp_2), \\
S_m &:= S_m^{(1)} \otimes I_K, \text{ with } S_m^{(1)} = \begin{bmatrix} I_m \\ O_{(p-m) \times m} \end{bmatrix} & (Kp \times Km), \\
\Pi'_m &:= \Pi'^{(1)}_m \otimes I_K, \text{ with } \Pi'^{(1)}_m = \begin{bmatrix} I_{m-p_1} \\ O_{(p-m) \times (m-p_1)} \end{bmatrix} & (Kp_2 \times K(m-p_1)).
\end{aligned}
\quad (2.3)
$$

The Kronecker products with $I_K$ are a direct consequence of estimating all equations with the same lag order. The regressor matrix for the estimation of a VAR($m$) model satisfies $Z_{T,m} = S_m' Z_T$. The implied OLS estimator is the ($K \times Km$) matrix $\hat{\Theta}_{T,m}$ given by

$$\hat{\Theta}_{T,m} = \Theta_{T,m} + C_T(I_{Kp_2} - \Pi_m' \Pi_m) S_0' Z_T Z_{T,m}' (Z_{T,m} Z_{T,m}')^{-1} + U Z_{T,m}' (Z_{T,m} Z_{T,m}')^{-1}. \tag{2.4}$$

Some rearranging and rescaling produces,

$$\sqrt{T}\left(\hat{\Theta}_{T,m} - \Theta_{T,m}\right) = \underbrace{\sqrt{T} C_T}_{\Delta} \left(I_{Kp_2} - \Pi_m' \Pi_m\right) S_0' \left(\frac{1}{T} Z_T Z_T'\right) S_m \left[S_m' \left(\frac{1}{T} Z_T Z_T'\right) S_m\right]^{-1}$$

$$+ \left(\frac{1}{\sqrt{T}} U Z_T'\right) S_m \left[S_m' \left(\frac{1}{T} Z_T Z_T'\right) S_m\right]^{-1}, \tag{2.5}$$

and it can be seen that the $T^{1/2}$-consistency of the estimator precisely matches the decay rate of $T^{-1/2}$ in the elements of the parameter matrix $C_T$. As a final step we apply the vec operator to transform the parameter matrices into a single parameter vector,

$$\sqrt{T}\left(\hat{\theta}_{T,m} - \theta_{T,m}\right) = \left(\left[S_m' \left(\frac{1}{T} Z_T Z_T'\right) S_m\right]^{-1} S_m' \left(\frac{1}{T} Z_T Z_T'\right) S_0 \left(I_{Kp_2} - \Pi_m' \Pi_m\right) \otimes I_K\right) \delta$$

$$+ \left(\left[S_m' \left(\frac{1}{T} Z_T Z_T'\right) S_m\right]^{-1} S_m' \otimes I_K\right) \frac{1}{\sqrt{T}} \sum_{t=1}^{T} \text{vec}\left(u_t z_{T,t-1}'\right), \tag{2.6}$$

where $\hat{\theta}_{T,m} = \text{vec}(\hat{\Theta}_{T,m})$, $\theta_{T,m} = \text{vec}(\Theta_{T,m})$, and $\delta = \text{vec}(\Delta)$. Equation (2.6) depends on: (1) various selection matrices, (2) the random matrix $\frac{1}{T} Z_T Z_T' = \frac{1}{T} \sum_{t=1}^{T} z_{T,t-1} z_{T,t-1}'$, and (3) the random vector $\frac{1}{\sqrt{T}} \sum_{t=1}^{T} \text{vec}(u_t z_{T,t-1}')$. The latter two rescaled sums are typically encountered in laws of large numbers and central limit theorems, respectively. The following three assumptions guarantee that such theorems are applicable.

**Assumption 1.** *The sequence $\{u_t\}$ of random K-vectors is an independent and identically distributed sequence with mean vector zero, a positive definite covariance matrix $E(u_t u_t') = \Sigma$, and there exists a $c > 0$ such that $E|u_{it} u_{jt} u_{kt} u_{mt}| < c < \infty$ for $i, j, k, m = 1, 2, \ldots K$.*

**Assumption 2.** $\det(B_T(z)) = \det\left(I_K z^p - B_1 z^{p-1} - \ldots - B_{p_1} z^{p_2} - \frac{\Delta_1}{\sqrt{T}} z^{p_2-1} - \ldots - \frac{\Delta_{p_2}}{\sqrt{T}}\right) \neq 0$ for all $|z| \geq 1$ and $\forall T \in \mathbb{N}$.

**Assumption 3.** $\det(B_\infty(z)) = \det\left(I_K z^{p_1} - B_1 z^{p_1-1} - \ldots - B_{p_1}\right) \neq 0$ for all $|z| \geq 1$.

Assumption 1 provides moment bounds and independence between the innovation $u_t$ and its past. This latter property is exploited to apply limit theorems for martingale differences.[2] Assumptions 2 and 3 require the vector autoregressive process to be stationary for every finite $T$ and also in the absence of local misspecification. The asymptotic properties of the OLS estimators are stated in Theorem 1.

---

[2]The requirement of i.i.d. innovations can be relaxed to the assumption that $\{u_t\}$ is a martingale difference sequence. Formally, let $\mathcal{F}_t = \sigma(u_s, -\infty < s \leq t)$ denote the sigma field generated by the innovations up to and including time $t$. Our results remain valid if the conditions $E(u_t) = 0$ and $E(u_t u_t') = \Sigma$ are replaced by $E(u_t|\mathcal{F}_{t-1}) = 0$ and $E(u_t u_t'|\mathcal{F}_{t-1}) = \Sigma$, respectively.

**Theorem 1** (Asymptotic Normality of the Least Squares Estimator). *Let Assumptions 1–3 hold. Then*
(a) *In the limit $T \to \infty$, we have for any $m \in \mathcal{M} = \{p_1, p_1 + 1, \ldots, p\}$,*

$$\sqrt{T} \left( \hat{\boldsymbol{\theta}}_{T,m} - \boldsymbol{\theta}_{T,m} \right) \overset{d}{\longrightarrow} A_m \boldsymbol{\delta} + \left( \left[ S'_m \boldsymbol{\Omega} S_m \right]^{-1} S'_m \otimes I_K \right) R \sim \mathrm{N} \left( A_m \boldsymbol{\delta}, \left[ S'_m \boldsymbol{\Omega} S_m \right]^{-1} \otimes \boldsymbol{\Sigma} \right),$$

*with $\boldsymbol{\Omega} = \mathrm{plim}_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} z_{T,t-1} z'_{T,t-1}$, $A_m = \left[ S'_m \boldsymbol{\Omega} S_m \right]^{-1} S'_m \boldsymbol{\Omega} S_0 \left( I_{Kp_2} - \boldsymbol{\Pi}'_m \boldsymbol{\Pi}_m \right) \otimes I_K$ and $R \sim \mathrm{N} \left( \mathbf{0}, \boldsymbol{\Omega} \otimes \boldsymbol{\Sigma} \right)$.*

(b) *Let $\hat{\boldsymbol{u}}^m_{T,t}$ denote the OLS residuals from the estimation of a VAR(m). Consider $\hat{\boldsymbol{\Sigma}}^m_u = \frac{1}{T} \sum_{t=1}^{T} \hat{\boldsymbol{u}}^m_{T,t} \hat{\boldsymbol{u}}^{m'}_{T,t}$ as an estimator for $\boldsymbol{\Sigma}_u$. The result in part (a) can be strengthened to joint asymptotic normality with the covariance matrix estimator $\hat{\boldsymbol{\Sigma}}^m_u$,*

$$\begin{bmatrix} \sqrt{T} \left( \hat{\boldsymbol{\theta}}_{T,m} - \boldsymbol{\theta}_{T,m} \right) \\ \sqrt{T} \mathrm{vech} \left( \hat{\boldsymbol{\Sigma}}^m_u - \boldsymbol{\Sigma} \right) \end{bmatrix} \overset{d}{\longrightarrow} \mathrm{N} \left( \begin{bmatrix} A_m \\ \mathbf{O} \end{bmatrix} \boldsymbol{\delta}, \begin{bmatrix} \left[ S'_m \boldsymbol{\Omega} S_m \right]^{-1} \otimes \boldsymbol{\Sigma} & \mathbf{O} \\ \mathbf{O} & \boldsymbol{\Xi}_{22} \end{bmatrix} \right).$$

*The matrix $\boldsymbol{\Xi}_{22}$ is specified in the Appendix.*

(c) *The estimator convergence as discussed in parts (a) and (b) of this theorem is also a joint convergence across different $m \in \mathcal{M}$. That is, for $\{i_1, i_2, \ldots, i_M\} \in \mathcal{M}$, any $m \in \mathcal{M}$, and $i_1 < i_2 < \ldots < i_M$, we have*

$$\begin{bmatrix} \sqrt{T} \left( \hat{\boldsymbol{\theta}}_{T,i_1} - \boldsymbol{\theta}_{i_1} \right) \\ \sqrt{T} \left( \hat{\boldsymbol{\theta}}_{T,i_2} - \boldsymbol{\theta}_{i_2} \right) \\ \vdots \\ \sqrt{T} \left( \hat{\boldsymbol{\theta}}_{T,i_M} - \boldsymbol{\theta}_{i_M} \right) \\ \sqrt{T} \mathrm{vech} \left( \hat{\boldsymbol{\Sigma}}^m_u - \boldsymbol{\Sigma} \right) \end{bmatrix} \overset{d}{\longrightarrow} \mathrm{N} \left( \begin{bmatrix} A_{i_1} \\ A_{i_2} \\ \vdots \\ A_{i_M} \\ \mathbf{O} \end{bmatrix} \boldsymbol{\delta}, \begin{bmatrix} V_{i_1 i_1} & V_{i_1 i_2} & \ldots & V_{i_1 i_M} & \mathbf{O} \\ V_{i_2 i_1} & V_{i_2 i_2} & \ldots & V_{i_2 i_M} & \mathbf{O} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ V_{i_M i_1} & V_{i_M i_2} & \ldots & V_{i_M i_M} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \ldots & \mathbf{O} & \boldsymbol{\Xi}_{22} \end{bmatrix} \right).$$

*The matrices $V_{jk}$ are given by $V_{jk} = \left[ S'_j \boldsymbol{\Omega} S_j \right]^{-1} S'_j \boldsymbol{\Omega} S_k \left[ S'_k \boldsymbol{\Omega} S_k \right]^{-1} \otimes \boldsymbol{\Sigma}$. It suffices to consider a single estimator for $\boldsymbol{\Sigma}$ because all the estimators are asymptotically equivalent.*

The matrix $\boldsymbol{\Omega}$ deserves further attention. It is defined as the probability limit of the Gram matrix $\frac{1}{T} Z_T Z'_T$. The proof of Theorem 1 reveals that this probability limit equals $\mathrm{E}(z_{\infty,t} z'_{\infty,t})$ where $z_{\infty,t}$ is defined as in Equation (2.2) but being generated by a VAR without local misspecification. We illustrate this remark with the AR process defined by $y_{T,t} = \alpha y_{T,t-1} + \frac{\delta_1}{\sqrt{T}} y_{T,t-2} + \frac{\delta_2}{\sqrt{T}} y_{T,t-3} + u_t$, that is an AR model with $p_1 = 1$, $p_2 = 2$ and $p = 3$: in that case $\boldsymbol{\Omega} = \frac{\sigma^2}{1-\alpha^2} \begin{bmatrix} 1 & \alpha & \alpha^2 \\ \alpha & 1 & \alpha \\ \alpha^2 & \alpha & 1 \end{bmatrix}$.

**Remark 1.** The consequences of the local misspecification framework are visible in Theorem 1. Standard asymptotics will fail if relevant parameters are left out since omitted variable bias will dominate asymptotically.[3] The local-to-zero rate of $T^{-1/2}$ balances this diverging behavior such that a finite asymptotic bias remains. This reasoning applies to all models that contain all the fixed parameters (i.e., the lag order should be no less than $p_1$) and leave out arbitrary parameters that are local-to-zero.

**Remark 2.** Assumption 2 is rather strict because it requires stationarity for all $T$ in the natural numbers. Is it even possible for any parameter combination to satisfy this assumption? We can answer this question

---

[3] Let us consider the data generating process $y_t = \alpha_1 y_{t-1} + \alpha_2 y_{t-2} + u_t$. Suppose that we estimate an AR(1) model. The OLS parameter estimator of the first lag coefficient, say $\hat{\alpha}$, satisfies $\sqrt{T}(\hat{\alpha} - \alpha_1) = \sqrt{T} \alpha_2 \frac{\frac{1}{T} \sum_{t=1}^{T} y_{t-1} y_{t-2}}{\frac{1}{T} \sum_{t=1}^{T} y_{t-1}^2} + \frac{\frac{1}{\sqrt{T}} \sum_{t=1}^{T} y_{t-1} u_t}{\frac{1}{T} \sum_{t=1}^{T} y_{t-1}^2}$. The first term on the RHS diverges at large $T$ for $\alpha_2 \neq 0$. The divergence rate is $\sqrt{T}$.

in the affirmative for the univariate case but the result does not generalize easily to the multivariate case. For the univariate case, we define the lag polynomial $\beta_T(L)$ by

$$\beta_T(L)y_{T,t} = \left(1 - \beta_1 L - \ldots - \beta_{p_1} L^{p_1} - \frac{\delta_1}{\sqrt{T}}L^{p_1+1} - \ldots - \frac{\delta_{p_2}}{\sqrt{T}}L^p\right)y_{T,t} = u_t. \tag{2.7}$$

Fujiwara (1916) has shown that the largest modulus root of a polynomial $a(z) = a_0 z^n + a_1 z^{n-1} + \ldots + a_{n-1}z + a_n$ is bounded above by $2\max\left\{|a_1/a_0|, |a_2/a_0|^{1/2}, \ldots, |a_n/a_0|^{1/n}\right\}$. The largest modulus root of the lag polynomial $\beta_T$ is thus bounded by

$$2\max\left\{|\beta_1|, |\beta_2|^{\frac{1}{2}}, \ldots, |\beta_{p_1}|^{\frac{1}{p_1}}, |\frac{\delta_1}{\sqrt{T}}|^{\frac{1}{p_1+1}}, \ldots, |\frac{\delta_{p_2}}{\sqrt{T}}|^{\frac{1}{p}}\right\}. \tag{2.8}$$

We deduce from Equation (2.8) that $2\max\left\{|\beta_1|, |\beta_2|^{\frac{1}{2}}, \ldots, |\beta_{p_1}|^{\frac{1}{p_1}}, |\delta_1|^{\frac{1}{p_1+1}}, \ldots, |\delta_{p_2}|^{\frac{1}{p}}\right\} < 1$ guarantees stationarity for all $T$.[4]

## 2.3. Quantities of interest

The focused information criterion (FIC) introduced by Claeskens and Hjort (2003) focusses on a quantity of interest rather than general model fit. Quantities of interest could be a single parameter, several parameters, or parameter transformations. Natural quantities of interest in the current dynamical setting are the impulse responses. In general, let $\boldsymbol{\mu} : \mathbb{R}^{K^2 p + K(K+1)/2} \to \mathbb{R}^l$ define the mapping from the model parameters to the $l$-dimensional focus quantity. The first $K^2 p$ arguments of the function $\boldsymbol{\mu}$ are reserved for the conditional mean parameters, whereas the last $K(K+1)/2$ arguments refer to the parameters in $\boldsymbol{\Sigma}$. As such we define $\boldsymbol{\sigma} = \text{vech}(\boldsymbol{\Sigma})$ and $\widehat{\boldsymbol{\sigma}} = \text{vech}(\widehat{\boldsymbol{\Sigma}})$, and write $\mu(\boldsymbol{\theta}, \boldsymbol{\sigma})$.[5] We additionally assume that evaluating the quantity of interest at $\boldsymbol{\mu}((\boldsymbol{\theta}_{T,m}, \mathbf{0}_{K^2(p-m)}), \boldsymbol{\sigma})$ provides an estimate for the quantity of interest in the model with $m$ lags. The Auxiliary Result in the Appendix shows that this is true for the impulse responses. The next theorem follows from Theorem 1 and the multivariate first-order delta method.

**Theorem 2** (Asymptotic Normality of the Quantities of Interest). *Let $\boldsymbol{\mu} : \mathbb{R}^{K^2 p + K(K+1)/2} \to \mathbb{R}^l$ have a continuous first derivative at all points $(\boldsymbol{\theta}_m, \mathbf{0}_{K^2(p-m)}, \boldsymbol{\sigma})$, with $m \in \mathcal{M}$. Let $\boldsymbol{\theta}_\infty$ denote the parameters obtained by taking $\boldsymbol{\theta}_{T,p}$ but setting $\boldsymbol{\Delta}_1 = \boldsymbol{\Delta}_2 = \ldots = \boldsymbol{\Delta}_{p_2} = \mathbf{O}$, and define the Jacobian matrices $\boldsymbol{D}_\theta = \partial\boldsymbol{\mu}(\boldsymbol{\theta}_\infty, \boldsymbol{\sigma})/\partial\boldsymbol{\theta}'$ and $\boldsymbol{D}_\sigma = \partial\boldsymbol{\mu}(\boldsymbol{\theta}_\infty, \boldsymbol{\sigma})/\partial\boldsymbol{\sigma}'$. For $\boldsymbol{D}_\theta$ and $\boldsymbol{D}_\sigma$ not having zero rows, under Assumptions 1–3, and as $T \to \infty$,*

$$\sqrt{T}\left(\boldsymbol{\mu}\left((\widehat{\boldsymbol{\theta}}_{T,m}, \mathbf{0}_{K^2(p-m)}), \widehat{\boldsymbol{\sigma}}\right) - \boldsymbol{\mu}(\boldsymbol{\theta}_{T,p}, \boldsymbol{\sigma})\right) \xrightarrow{d} \text{N}\left(\boldsymbol{D}_\theta \boldsymbol{C}_m \boldsymbol{\delta}, \boldsymbol{D}_\theta \boldsymbol{P}_m(\boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})\boldsymbol{P}_m \boldsymbol{D}_\theta' + \boldsymbol{D}_\sigma \boldsymbol{\Xi}_{22} \boldsymbol{D}_\sigma'\right),$$

*with $\boldsymbol{P}_m = \boldsymbol{S}_m\left[\boldsymbol{S}_m'\boldsymbol{\Omega}\boldsymbol{S}_m\right]^{-1}\boldsymbol{S}_m' \otimes \boldsymbol{I}_K$ and*

$$\boldsymbol{C}_m = \left(\boldsymbol{S}_m\left[\boldsymbol{S}_m'\boldsymbol{\Omega}\boldsymbol{S}_m\right]^{-1}\boldsymbol{S}_m'\boldsymbol{\Omega} - \boldsymbol{I}_{Kp}\right)\boldsymbol{S}_0\left(\boldsymbol{I}_{Kp_2} - \boldsymbol{\Pi}_m'\boldsymbol{\Pi}_m\right) \otimes \boldsymbol{I}_K.$$

We define the impulse response at horizon $h$ as the $h$th coefficient matrix of the MA($\infty$) representation $\boldsymbol{y}_t = \sum_{i=0}^{\infty} \boldsymbol{\Phi}_i \boldsymbol{u}_{t-i}$ with $\boldsymbol{\Phi}_0 = \boldsymbol{I}_K$, hence $\boldsymbol{\mu} : \mathbb{R}^{K^2 p + K(K+1)/2} \to \mathbb{R}^{K^2}$. Theorem 2 can be applied if the Jacobian matrices $\boldsymbol{D}_\theta$ and $\boldsymbol{D}_\sigma$ are known. Lütkepohl (1990) lists these Jacobian matrices for the impulse responses, the orthogonalized impulse responses, the accumulated responses, the total accumulated

---

[4]This condition is a sufficient but by no means a necessary condition. For $p_1 = 1$ and $p_2 = 1$, the model $y_{T,t} = 0.7y_{T,t-1} + \frac{0.75}{\sqrt{T}}y_{T,t-2} + u_t$ is stationary for all $T$ but the parameters violate the requirement based on Fujiwara's bound.

[5]Theorem 1 showed that all the $\widehat{\boldsymbol{\Sigma}}_u^m$ are asymptotically equivalent. We omit the superscript $m$ from now on.

responses, and the forecast error variance decomposition. The specific case of the (orthogonalized) impulse responses is highlighted in the following Corollary.

**Corollary** (An Application to Impulse Responses). *Let $A_\infty$ denote the $(Kp \times Kp)$ companion matrix related to the process in which the misspecification coefficients have been set to zero, $J = [I_K \; O \; \cdots \; O]$ a matrix of dimensions $(K \times Kp)$ and $\Sigma_u = PP'$. Then, under the assumptions of Theorem 2,*
(a) *The asymptotic distribution of the estimated impulse response at horizon $i$, $\hat{\Phi}_i$, follows*

$$\sqrt{T}\left(\text{vec}(\hat{\Phi}_i) - \text{vec}(\Phi_i)\right) \xrightarrow{d} \text{N}\left(G_i C_m \delta, G_i P_m \left(\Omega \otimes \Sigma\right) P_m G_i'\right),$$

*where $G_i = \partial\text{vec}\left(\Phi_i\right)/\partial\theta = \sum_{j=0}^{i-1} J\left(A_\infty\right)^{i-1-j} \otimes \Phi_j$.*
(b) *The asymptotic distribution of the estimated orthogonalized impulse response at horizon $i$, $\hat{\Psi}_i$, follows*

$$\sqrt{T}\left(\text{vec}(\hat{\Psi}_i) - \text{vec}(\Psi_i)\right) \xrightarrow{d} \text{N}\left(F_i C_m \delta, F_i P_m \left(\Omega \otimes \Sigma\right) P_m F_i' + \bar{F}_i \Xi_{22} \bar{F}_i'\right),$$

*where $F_0 = O$ and $F_i = \left(P' \otimes I_K\right) G_i$ for $i > 0$. For all $i$ we have $\bar{F}_i = \left(I_K \otimes \Phi_i\right) H$ with $H = \partial\text{vec}(P)/\partial\sigma' = L_K'[L_K\left(I_{K^2} + K_{KK}\right)\left(P \otimes I_K\right) L_K']^{-1}$ (see Lütkepohl (1990) for the definitions of $L_K$ and $K_{KK}$).*

**Remark 3.** The first-order delta method is invalid if either $D_\theta$ or $D_\sigma$ has zero rows. It is well-documented in the literature that $D_\theta$ can have zero rows for specific parameter combinations when impulse responses are considered. We refer to Lütkepohl (1990) and Benkwitz et al. (2000) for details.

### *2.3.1. Model selection: The focused information criterion (FIC)*
The intuition behind the FIC of Claeskens and Hjort (2003) is most easily understood for a univariate quantity of interest, so we temporarily assume $l = 1$. The generalization to multiple quantities is covered in Remark 4. Theorem 2 implies that the asymptotic mean squared error (AMSE) of the focus quantity $\mu\left((\hat{\theta}_{T,m}, \mathbf{0}_{K^2(p-m)}), \hat{\sigma}\right)$ is

$$\text{AMSE}\left(\mu\left((\hat{\theta}_{T,m}, \mathbf{0}_{K^2(p-m)}), \hat{\sigma}\right)\right) = D_\theta \left[C_m \delta\delta' C_m' + P_m(\Omega \otimes \Sigma)P_m\right] D_\theta' + D_\sigma \Xi_{22} D_\sigma'. \tag{2.9}$$

There are three contributions to the AMSE: (1) the term $D_\theta C_m \delta\delta' C_m' D_\theta'$ is an asymptotic squared bias originating from the exclusion of local-to-zero parameters, (2) the asymptotic variance contribution $D_\theta P_m(\Omega \otimes \Sigma)P_m D_\theta'$, and (3) the contribution $D_\sigma \Xi_{22} D_\sigma'$ which does not depend on the lag order $m$. Overall we face a bias-variance tradeoff when having to decide on $m$.

The FIC is an estimate of the AMSE. The quantities $\hat{\theta}_{T,p}$ and $\hat{\Omega} = \frac{1}{T}\sum_{t=1}^T z_{T,t} z_{T,t}'$ provide consistent estimates for $\theta_T$ and $\Omega$, respectively. In view of the continuous mapping theorem, $\hat{P}_m = S_m \left[S_m' \hat{\Omega} S_m\right]^{-1} S_m' \otimes I_K$ and

$$\hat{C}_m = \left(S_m \left[S_m' \hat{\Omega} S_m\right]^{-1} S_m' \hat{\Omega} - I_{Kp}\right) S_0 \left(I_{Kp_2} - \Pi_m' \Pi_m\right) \otimes I_K, \tag{2.10}$$

are consistent estimators as well. *A consistent estimator for $\delta$ is not available due to the adopted misspecification framework.* We follow the existing literature (see Claeskens and Hjort (2003), Liu (2015), and Charkhi et al. (2016) among others) and use $\hat{\delta} = \sqrt{T}\text{vec}(\hat{\Theta}_{T,p} S_0)$ which satisfies[6]

$$\hat{\delta} \xrightarrow{d} R_\delta = \delta + (S_0' \Omega^{-1} \otimes I_K)R \sim \text{N}(\delta, S_0' \Omega^{-1} S_0 \otimes \Sigma). \tag{2.11}$$

---

[6]$\hat{\delta} = \sqrt{T}\,\text{vec}(\hat{\Theta}_{T,p} S_0)$ is the sample equivalent of $\delta = \text{vec}(\Delta) = \sqrt{T}\text{vec}(C_T) = \sqrt{T}\text{vec}(\Theta_{T,p} S_0)$.

Asymptotically, we have $E(\hat{\hat{\delta}}\hat{\hat{\delta}}') = \delta\delta' + S_0'\Omega^{-1}S_0 \otimes \Sigma$. Using the asymptotically unbiased estimate $\hat{\hat{\delta}}\hat{\hat{\delta}}' - S_0'\hat{\Omega}^{-1}S_0 \otimes \hat{\Sigma}$ for $\delta\delta'$, the FIC for the approximating model with $m$ lags is defined as

$$\widehat{FIC}_m = D_\theta \left[ \hat{C}_m \left( \hat{\hat{\delta}}\hat{\hat{\delta}}' - S_0'\hat{\Omega}^{-1}S_0 \otimes \hat{\Sigma} \right) \hat{C}_m' + \hat{P}_m(\hat{\Omega} \otimes \hat{\Sigma})\hat{P}_m \right] D_\theta' + D_\sigma \hat{\Xi}_{22}D_\sigma'. \tag{2.12}$$

This estimate of the AMSE can be computed for every model, and the model with the lowest $\widehat{FIC}_m$ is selected. We elaborate on implications of inconsistent estimation of $\delta$ in Section 2.4.

**Remark 4.** The same procedure can be followed when $l > 1$, but the AMSE becomes an $(l \times l)$ matrix. The trace or determinant are meaningful ways to describe this AMSE matrix by a see scalar Charkhi et al. (2016).[7] The trace is computationally convenient because the overall FIC will be the sum of the individual univariate FIC contributions.

### 2.3.2. Model averaging: Plug-in averaging

Liu (2015) proposed a model averaging approach along the lines of the FIC. It was named plug-in averaging. We again depart from the case $l = 1$, see Remark 5 for the generalization. Part (b) from Theorem 1 implies that linear combinations of the VAR parameter estimates are also asymptotically normally distributed. Interpret the coefficients in the linear combination as weights, i.e., define $w = (w_{p_1}, w_{p_1+1}, \ldots, w_p)$ with $w \in \mathcal{H} = \left\{ w \in [0,1]^{p_2+1} : \sum_{m=p_1}^{p} w_m = 1 \right\}$.[8]

Theorem 3 details the asymptotic distribution of the following weighted estimator,

$$\bar{\mu}(w) = \sum_{m=p_1}^{p} w_m \mu\big((\hat{\theta}_{T,m}, \mathbf{0}_{K^2(p-m)}), \hat{\sigma}\big). \tag{2.13}$$

**Theorem 3** (Asymptotic Normality of the Plug-In Estimator). *Under the Assumptions of Theorem 2, we have for $T \to \infty$,*

$$\sqrt{T}\big(\bar{\mu}(w) - \mu(\theta_T, \sigma)\big) \xrightarrow{d} D_\theta \sum_{m=p_1}^{p} w_m C_m \delta + D_\theta \sum_{m=p_1}^{p} w_m P_m R + D_\sigma S$$

$$\sim N\left( D_\theta \sum_{m=p_1}^{p} w_m C_m \delta, V + D_\sigma \Xi_{22}D_\sigma' \right),$$

*with $V = \sum_{m=p_1}^{p} \sum_{l=p_1}^{p} w_m w_l D_\theta P_m (\Omega \otimes \Sigma) P_l D_\theta'$.*

As for the FIC, we compute the AMSE and find[9]

$$\text{AMSE}(\bar{\mu}(w)) = \sum_{m=p_1}^{p} \sum_{l=p_1}^{p} w_m D_\theta \big( C_m \delta\delta' C_l' + P_m (\Omega \otimes \Sigma) P_l \big) D_\theta' w_l = w'\Psi w, \tag{2.14}$$

with the $((p_2 + 1) \times (p_2 + 1))$ matrix $\Psi$ having the $(m,l)$th element $\Psi_{m,l} = D_\theta \big( C_m \delta\delta' C_l' + P_m (\Omega \otimes \Sigma) P_l \big) D_\theta'$. The optimal weight vector that minimizes the AMSE is given by $w^0 = \arg\min_{w \in \mathcal{H}} w'\Psi w$. But $w^0$ depends on population quantities, so using the same estimates for population

---

[7] Any mapping from the AMSE matrix to a scalar can be used, e.g., matrix norms could be used as well.

[8] We assume that we average over all the models in $\mathcal{M} = \{p_1, p_1 + 1, \ldots, p\}$.

[9] The contribution $D_\sigma \Xi_{22}D_\sigma'$ does not depend on $m$, and is therefore inconsequential for the analysis. This term will be omitted from the AMSE to allow for an easier presentation.

quantities as in Section 2.3.1, we compute feasible weights as

$$\hat{w} = \arg\min_{w \in \mathcal{H}} w' \hat{\Psi}^j w, \qquad j \in \{\text{Biased, Bias cor}\}, \tag{2.15}$$

with

$$
\begin{aligned}
\hat{\Psi}_{m,l}^{\text{Biased}} &= D_\theta \left[ \hat{C}_m \hat{\delta}\hat{\delta}' \hat{C}_l' + \hat{P}_m \left( \hat{\Omega} \otimes \hat{\Sigma} \right) \hat{P}_l \right] D_\theta', \\
\hat{\Psi}_{m,l}^{\text{Bias cor}} &= D_\theta \left[ \hat{C}_m \left( \hat{\delta}\hat{\delta}' - S_0' \hat{\Omega}^{-1} S_0 \otimes \hat{\Sigma} \right) \hat{C}_l' + \hat{P}_m \left( \hat{\Omega} \otimes \hat{\Sigma} \right) \hat{P}_l \right] D_\theta'.
\end{aligned}
\tag{2.16}
$$

The sole difference between the matrix elements in Equation (2.16) is an asymptotic bias correction for $\hat{\delta}\hat{\delta}'$. For both versions, Equation (2.15) is a quadratic programming problem with linear constraints. Solvers are readily available (for example "quadprog" in Matlab, or "qprog" in Gauss). The estimator for $\delta$ remains inconsistent and we again refer to Section 2.4 for a discussion of the implications.

**Remark 5.** As in Remark 4, we will obtain an $(l \times l)$ AMSE matrix for multiple quantities of interest. This matrix has to be summarized by a scalar. The trace has again computational benefits because the objective function will take the form $\hat{w} = \arg\min_{w \in \mathcal{H}} w' \left( \sum_i \hat{\Psi}_i^j \right) w$ with $\hat{\Psi}_i^j$ the matrix corresponding to the $i$th focus quantity.

**Remark 6.** Two remarks related to Charkhi et al. (2016) are in place.
1. The weight vector $\hat{w}$ is only uniquely determined when $\hat{\Psi}^j$ is positive definite. As such, the bias subtraction may lead to nonunique weights.
2. Charkhi et al. (2016) consider a weighting scheme in which the weights sum to one but are not necessarily positive. Simulation results have shown that it is advisable to keep the positivity constraint in our autoregressive setup because weights can otherwise become large in magnitude and rather unstable.

**Remark 7.** Autoregressive models of infinite order have been considered by Berk (1974) and Lewis and Reinsel (1985) among others. It is an intriguing question whether the current framework can be extended to VAR($\infty$) models.[10] We argue that the main difficulty is the estimation of the infinitely many local-to-zero parameters. Let us consider the univariate model $y_{T,t} = \alpha y_{T,t-1} + \sum_{j=1}^{\infty} \left( \frac{\delta_j}{\sqrt{T}} \right) y_{T,t-1-j} + u_t$ as an illustration. We conjecture[11] that the asymptotic distribution of the approximating AR(1) model follows $\sqrt{T}(\hat{\alpha} - \alpha) \xrightarrow{d} N \left( \sum_{j=1}^{\infty} \delta_j \alpha^j, 1 - \alpha^2 \right)$. The bias contribution to the AMSE now depends on infinitely many $\delta_j$. Their estimation would require the lag order of the largest approximating model to grow with sample size. Our proof of Theorem 1 does not easily allow for such an extension since we currently rely on the finite dimension of the companion matrix. A full exploration of this topic is left for further research.

There is one final result that forces us to look at the case $l > 1$. Practitioners are usually interested in the impulse responses for several horizons. Using a separate weight vector for every horizon may: (1) create impulse responses that vary irregularly from one horizon to the next due to strong changes in the weights, and (2) result in confidence intervals that do not take into account the dependence between the horizons. Theorem 4 extends the result of Liu (2015) to obtain asymptotically valid confidence bands for several horizons.

---

[10]We thank an anonymous referee for bringing this topic to our attention.
[11]We can be more precise concerning our assumptions. Theorem 1 has shown that the asymptotic results are governed by the process with the local-to-zero parameters being set equal to zero. We assume that this remains true when there are infinitely many local-to-zero parameters.

**Theorem 4** (Joint Confidence Bands). *Under the Assumptions of Theorem 2, if $w_m(\hat{\boldsymbol{\delta}}) \xrightarrow{d} w_m(\boldsymbol{R}_\delta)$, and if $\boldsymbol{D}_\theta\,(\boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})\,\boldsymbol{D}_\theta' + \boldsymbol{D}_\sigma\,\boldsymbol{\Xi}_{22}\boldsymbol{D}_\sigma \succ 0$, then*

$$\left(\sqrt{T}\Big(\bar{\boldsymbol{\mu}}(\hat{\boldsymbol{w}}) - \boldsymbol{\mu}(\boldsymbol{\theta}_T, \boldsymbol{\sigma})\Big) - \boldsymbol{D}_\theta \sum_{m=p_1}^{p} \hat{w}_m \hat{\boldsymbol{C}}_m \hat{\boldsymbol{\delta}}\right)' \left(\boldsymbol{D}_\theta \left(\hat{\boldsymbol{\Omega}} \otimes \hat{\boldsymbol{\Sigma}}\right) \boldsymbol{D}_\theta' + \boldsymbol{D}_\sigma\,\hat{\boldsymbol{\Xi}}_{22}\boldsymbol{D}_\sigma\right)^{-1}$$

$$\left(\sqrt{T}\Big(\bar{\boldsymbol{\mu}}(\hat{\boldsymbol{w}}) - \boldsymbol{\mu}(\boldsymbol{\theta}_T, \boldsymbol{\sigma})\Big) - \boldsymbol{D}_\theta \sum_{m=p_1}^{p} \hat{w}_m \hat{\boldsymbol{C}}_m \hat{\boldsymbol{\delta}}\right) \leq \chi^2_{l, 1-\alpha},$$

*is an asymptotically correct $(1 - \alpha)$ confidence band, where $\chi^2_{l,1-\alpha}$ denotes the $(1 - \alpha)$ quantile of a chi-squared distributed random variable with l degrees of freedom.*

**Remark 8.** There is one practical concern which has not been addressed, namely the choices for $p_1$ and $p$.[12] $p_1$ will turn out to be unimportant. To see this, we consider the expression for $\widehat{FIC}_m$ (a similar reasoning applies to the plug-in weights). The terms $\boldsymbol{D}_\theta \hat{\boldsymbol{P}}_m(\hat{\boldsymbol{\Omega}} \otimes \hat{\boldsymbol{\Sigma}})\hat{\boldsymbol{P}}_m \boldsymbol{D}_\theta'$ and $\boldsymbol{D}_\sigma\,\hat{\boldsymbol{\Xi}}_{22}\boldsymbol{D}_\sigma'$ in Equation (2.12) do not depend on $p_1$, so it remains to inspect the contribution $\boldsymbol{D}_\theta \hat{\boldsymbol{C}}_m \left(\hat{\boldsymbol{\delta}}\hat{\boldsymbol{\delta}}' - \boldsymbol{S}_0'\hat{\boldsymbol{\Omega}}^{-1}\boldsymbol{S}_0 \otimes \hat{\boldsymbol{\Sigma}}\right) \hat{\boldsymbol{C}}_m' \boldsymbol{D}_\theta'$. Using $\hat{\boldsymbol{\delta}} = \sqrt{T}(\boldsymbol{S}_0' \otimes \boldsymbol{I}_K)\hat{\boldsymbol{\theta}}_{T,p}$ we can rewrite this contribution as

$$\boldsymbol{D}_\theta \left(\hat{\boldsymbol{C}}_m(\boldsymbol{S}_0' \otimes \boldsymbol{I}_K)\right) \left[(\sqrt{T}\hat{\boldsymbol{\theta}}_{T,p})(\sqrt{T}\hat{\boldsymbol{\theta}}_{T,p})' - \hat{\boldsymbol{\Omega}}^{-1} \otimes \hat{\boldsymbol{\Sigma}}\right] \left(\hat{\boldsymbol{C}}_m(\boldsymbol{S}_0' \otimes \boldsymbol{I}_K)\right)' \boldsymbol{D}_\theta'. \qquad (2.17)$$

By definition of $\hat{\boldsymbol{C}}_m$, we have

$$\hat{\boldsymbol{C}}_m(\boldsymbol{S}_0' \otimes \boldsymbol{I}_K) = \left(\boldsymbol{S}_m \left[\boldsymbol{S}_m'\hat{\boldsymbol{\Omega}}\boldsymbol{S}_m\right]^{-1} \boldsymbol{S}_m'\hat{\boldsymbol{\Omega}} - \boldsymbol{I}_{Kp}\right) \left[\boldsymbol{S}_0 \left(\boldsymbol{I}_{Kp_2} - \boldsymbol{\Pi}_m'\boldsymbol{\Pi}_m\right)\boldsymbol{S}_0'\right] \otimes \boldsymbol{I}_K$$

$$= \left(\boldsymbol{S}_m \left[\boldsymbol{S}_m'\hat{\boldsymbol{\Omega}}\boldsymbol{S}_m\right]^{-1} \boldsymbol{S}_m'\hat{\boldsymbol{\Omega}} - \boldsymbol{I}_{Kp}\right) \underbrace{\left(\begin{bmatrix} \boldsymbol{O} & \boldsymbol{O} \\ \boldsymbol{O} & \boldsymbol{I}_{K(p-m)} \end{bmatrix} \otimes \boldsymbol{I}_K\right)}_{\mathcal{B}}, \qquad (2.18)$$

thereby showing that actually none of the contributions to $\widehat{FIC}_m$ depends on $p_1$. However, the zero pattern of the matrix $\mathcal{B}$ in Equation (2.18) will only cause a nondiverging value for $\widehat{FIC}_m$ if models are chosen such that $m \in \mathcal{M} = \{p_1, p_1 + 1, \ldots, p\}$. This supports the claim in Remark 1. The lag order of the full model, that is $p$, might be chosen by AIC or set equal to an a priori selected $p_{max}$.

### 2.4. Effects of inconsistently estimating delta

Equation (2.11) showed that $\hat{\boldsymbol{\delta}}$ converges to a normally distributed random vector centered around $\boldsymbol{\delta}$. How does this influence the selection and averaging procedures? Clearly, $\widehat{FIC}_m$, $\hat{\Psi}_{m,l}^{\text{Biased}}$, and $\hat{\Psi}_{m,l}^{\text{Bias cor}}$ will not converge in probability to the AMSE they are intended to estimate. The limiting distributions of these quantities are highlighted in Theorems 5 and 6. The plug-in results are stated for $\hat{\Psi}_{m,l}^{\text{Bias cor}}$, but a simple omission of the bias correction term would give the corresponding findings for $\hat{\Psi}_{m,l}^{\text{Biased}}$.

**Theorem 5** ([The Asymptotic Behavior of $\widehat{FIC}_m$). *] Under the Assumptions of Theorem 2, we have for $m \in \mathcal{M}\backslash p$,*

$$\widehat{FIC}_m \xrightarrow{d} \boldsymbol{D}_\theta \left[\boldsymbol{C}_m \left(\boldsymbol{R}_\delta \boldsymbol{R}_\delta' - (\boldsymbol{S}_0'\boldsymbol{\Omega}\boldsymbol{S}_0 \otimes \boldsymbol{\Sigma})\right) \boldsymbol{C}_m' + \boldsymbol{P}_m(\boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})\boldsymbol{P}_m\right] \boldsymbol{D}_\theta' + \boldsymbol{D}_\sigma\,\boldsymbol{\Xi}_{22}\boldsymbol{D}_\sigma'$$

$$:= FIC_m^\infty \sim a_m \chi^2_{noncentral}\left(1, \left(\boldsymbol{D}_\theta \boldsymbol{C}_m \boldsymbol{\delta}\right)^2 / a_m\right) - a_m + \boldsymbol{D}_\theta \boldsymbol{P}_m(\boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})\boldsymbol{P}_m \boldsymbol{D}_\theta' + \boldsymbol{D}_\sigma\,\boldsymbol{\Xi}_{22}\boldsymbol{D}_\sigma',$$

---

[12]We thank an anonymous referee who rightfully conjectured that the choice of $p_1$ does not have an influence on the numerical outcome for $\widehat{FIC}_m$.

where $a_m = D_\theta C_m(S_0'\Omega^{-1}S_0 \otimes \Sigma)C_m'D_\theta'$, and $\chi^2_{noncentral}(\nu,\lambda)$ denotes a noncentral chi-squared distributed random variable with $\nu$ degrees of freedom and noncentrality parameter $\lambda$. It can be shown that $\mathrm{E}\left(FIC_m^\infty\right) = AMSE\left(\mu(\hat{\theta}_{T,m}, \mathbf{0}_{K^2(p-m)}, \hat{\sigma})\right)$ and $\mathrm{var}(FIC_m^\infty) = 2a_m(a_m + 2(D_\theta C_m\delta)^2)$. For the full model, $m = p$, we have

$$\widehat{FIC}_p \xrightarrow{p} FIC_p^\infty = D_\theta(\Omega^{-1} \otimes \Sigma)D_\theta' = AMSE\left(\mu(\hat{\theta}_{T,p}, \hat{\sigma})\right).$$

**Theorem 6** (The Asymptotic Behavior of $\hat{\Psi}_{m,l}^{\text{Biased}}$ and $\hat{\Psi}_{m,l}^{\text{Bias cor}}$). *Under the Assumptions of Theorem 2, we have for $m, l \in \mathcal{M} \backslash p$,*

$$\hat{\Psi}_{m,l}^{Bias cor} \xrightarrow{d} R_\delta'C_m'D_\theta'D_\theta C_l R_\delta + D_\theta P_m(\Omega \otimes \Sigma)P_l'D_\theta' - D_\theta C_m(S_0'\Omega S_0 \otimes \Sigma)C_m'D_\theta'$$
$$:= \Psi_{m,l}^{Bias cor,\infty}.$$

*Two cases can be distinguished:*

(a) *If $m = l$, then $\Psi_{m,m}^{Bias cor,\infty} \sim a_m\chi^2_{noncentral}\left(1, \left(D_\theta C_m\delta\right)^2 / a_m\right) - a_m + D_\theta P_m(\Omega \otimes \Sigma)P_m'D_\theta'$.*

(b) *Define $\mathcal{A} = S_0'\Omega^{-1}S_0 \otimes \Sigma$, and consider the eigenvalue decomposition $\frac{1}{2}\mathcal{A}^{1/2}(C_m'D_\theta'D_\theta C_l + C_l'D_\theta'D_\theta C_m)\mathcal{A}^{1/2} = \sum_{i=1}^2 \lambda_i v_i v_i'$, where $\lambda_i$ denotes the eigenvalue corresponding to the eigenvector $v_i$. If $m \neq l$, then*

$$\Psi_{m,l}^{Bias cor,\infty} \sim \sum_{i=1}^2 \lambda_i\chi^2_{noncentral}\left(1, \left(v_i'\mathcal{A}^{-1/2}\delta\right)^2\right) + D_\theta P_m(\Omega \otimes \Sigma)P_l'D_\theta'$$
$$- D_\theta C_m(S_0'\Omega S_0 \otimes \Sigma)C_m'D_\theta'.$$

*If $m = p$ and/or $l = p$, then $\hat{\Psi}_{m,l}^{Bias cor} \xrightarrow{p} \Psi_{m,l}$.*

Theorems 5 and 6 stated the limiting distribution of the FIC and the matrix elements that enter the weighting scheme. Based on the random limits of these quantities, we might expect that our methods will not truly minimize the AMSE among either model choices or model weights. We proceed with a small illustration to stress the difference between knowing $\delta$ and having an estimator $\hat{\delta}$ that converges in distribution only.

### An illustration

Consider a simplified DGP, $y_{T,t} = 0.5y_{T,t-1} + \frac{\delta}{\sqrt{T}}y_{T,t-2} + u_t$, with $\mathrm{var}(u_t) = 1$, and a focus on the impulse response at horizon 1 (i.e., $D_\theta = (1,0)$ and $D_\sigma = 0$). The model set is $\mathcal{M} = \{1,2\}$. This simplified setting makes the asymptotic behavior of the FIC and plug-in weights analytically tractable. Figure 1(a) depicts $FIC_1^\infty$ and $FIC_2^\infty$ as a function of $\delta$. Note that $FIC_1^\infty$ converges in distribution and has a nonzero probability to give an outcome below $FIC_2^\infty$. This asymptotic selection probability of the model with $m = 1$ can be calculated analytically using Theorem 5. Figure 1(b) shows that the FIC does not select the model with the smallest AMSE with probability one.

Our simplified model can also be used to examine the effect of $\hat{\delta}$ on the plug-in weights. We focus on the weights in the absence of bias correction.[13] The $(2 \times 2)$ limiting matrix $\Psi^\infty$ is

$$\Psi^\infty = \begin{bmatrix} a_1\chi^2_{noncentral}\left(1, (D_\theta C_1\delta)^2/a_1\right) + \sigma^2 D_\theta P_1 D_\theta' & \sigma^2 D_\theta P_1 D_\theta' \\ \sigma^2 D_\theta P_1 D_\theta' & \sigma^2 D_\theta \Omega^{-1} D_\theta' \end{bmatrix}, \quad (2.19)$$

---

[13] There is a finite probability for the matrix $\Phi^\infty$ to have a negative eigenvalue when the bias correction is applied. This severely complicates the derivations, so we exclude this case from our analysis. For the $\Psi^\infty$ matrix without bias correction we will have $\Psi^\infty \succ 0$ if $D_\theta'(\Omega^{-1} - P_1)D_\theta > 0$. The latter requirement is equivalent to $D_\theta'\begin{bmatrix} -\omega_{12} \\ \omega_{11} \end{bmatrix} \neq 0$.

**Figure 1.** (a) The asymptotic MSE of the models with one and two lags (red and black line, respectively). The area between the 5 and 95% quantiles of $FIC_1^\infty$ is shaded in red. $\widehat{FIC}_2$ converges in probability to the values of the black line. (b) The asymptotic selection probabilities of the model with $m = 1$. The infeasible estimator takes a binary decision based on whether the red or black line in graph (a) is lowest. Model selection based on the focused information criterion results in a smoothed asymptotic selection probability because $\widehat{FIC}_1$ converges in distribution.



**Figure 2.** The 5 and 95% quantiles of the asymptotic distribution of the weights as a shaded red area, see Equation (2.20). The solid cyan lines are the asymptotically optimal weights which can only be obtained if either $\delta$ is known (infeasible) or a consistent estimator for $\delta$ is available.

where $a_1 = \sigma^2 D_\theta C_1 S_0' \Omega S_0 C_1' D_\theta'$ (see Theorem 6). Let $w^*$ denote the asymptotically optimal plug-in weight for the model with $m = 1$. We have

$$\Pr\left(w^* \leq x\right) = \Pr\left(\chi^2_{noncentral}\left(1, (D_\theta C_1 \delta)^2 / a_1\right) \geq \frac{\sigma^2 D_\theta (\Omega^{-1} - P_1) D_\theta' [1 - x]}{a_1 x}\right). \tag{2.20}$$

Figure 2 shows the area between the 5 and 95% quantiles of $w^*$ together with the optimal weight for known delta. We see that the asymptotic distribution of $w^*$ is located closer to zero than the optimal infeasible weights. This is unsurprising because the lack of bias-correction causes (on average) an overestimation of the AMSE of the model with $m = 1$.

**Remark 9.** The exposition in this section was based on a simplified model. We concluded that the absence of a consistent estimator for $\delta$ translates into suboptimal model selection and suboptimal model averaging. Also in more elaborate models, the FIC and the elements of the weighting matrix $\Psi$ will converge to random variables (except for $m = p$). It is key to realize that the AMSE with estimated $\delta$ will not coincide with the AMSE that can be attained if $\delta$ was either known or consistently estimated. We conjecture that these considerations are equally relevant outside an autoregressive framework, e.g., in the regression framework discussed in Liu (2015) and the likelihood framework of Charkhi et al. (2016).

## 3. Simulations

This simulation section consists of three parts.[14] In the first part we will verify our derivations for the simplified DGP, and see how the suboptimal selection/averaging affects the finite sample MSE. This section is followed by a study of the impulse responses for different horizons in a univariate and multivariate setting. All our graphs are made as a function of the scalar $\delta$. This scalar measures the amount of misspecification and the closeness to unit root.[15] Any missing starting values in the autoregressive recursion were replaced by zeros, and the first 100 data points were omitted as a presample. All results are based on 100,000 Monte Carlo replications.

The performance of the various methods was assessed using the empirical mean squared error. For model selection the featured methods are:

1. The Akaike information criterion ("**AIC**") and Bayesian information criterion ("**BIC**"), e.g., Section 4.3 of Lütkepohl (2005) and the original papers by Akaike (1998) and Schwarz et al. (1978).
2. The "**FIC**" from Equation (2.12) with estimated $\delta$.
3. An infeasible version of the FIC abbreviated as "**Infeas**." This information criterion is based on population quantities, especially $\delta$ is known.

For the model averaging setup we consider:

1. "**sAIC**" and "**sBIC**" as smoothed counterparts of the AIC and BIC, see Burnham and Anderson (2002). To illustrate, let $AIC(m)$ denote the AIC for model $m \in \mathcal{M}$. The smoothed AIC assigns a weight of $\exp\left(-\frac{1}{2}AIC(m)\right) / \sum_{m \in \mathcal{M}} \exp\left(-\frac{1}{2}AIC(m)\right)$ to model $m$.
2. Three plug-in averages are reported. "**Plug-in**" and "**Plug-in Corr.**" are computed from Equations (2.15) and (2.16), where only the second average uses the bias correction on $\hat{\delta}\hat{\delta}'$. The plug-in average based on known $\delta$ is denoted "**Infeas**."
3. The "**Jackknife**" model averaging procedure detailed in Hansen and Racine (2012) and Zhang et al. (2013).
4. The Stein combination shrinkage method used in the simulation section of Hansen (2016) is abbreviated "**SteinH**." This shrinkage method combines VAR(1) through VAR($p$) models as well as univariate AR(1) through AR($p$) models. Our DGP contains considerable interaction between the cross-sectional units so we also consider a shrinkage method abbreviated "**Stein**" which only combines the VAR(1) through VAR($p$).

---

[14]A selection of simulation results is reported here, extensive results can be found in the Supplementary Material.

[15]Previous studies (e.g., Hansen (2007), Hansen (2008), Hansen and Racine (2012), Liu and Okui (2013), Zhang et al. (2013) and Liu (2015)) show the performance as a function of the population $R^2$. This representation is inconvenient in our dynamic setup because it is unclear when we are approaching the boundary of the stationarity region.
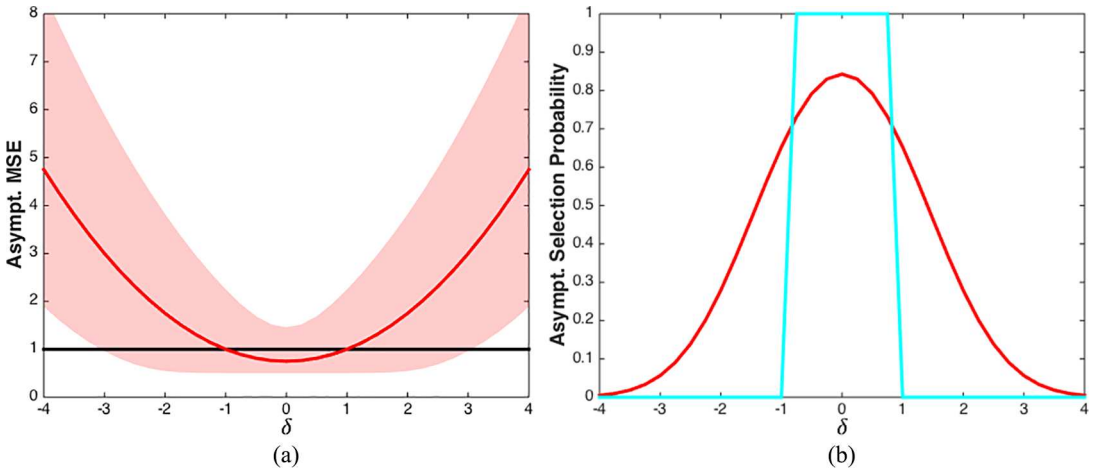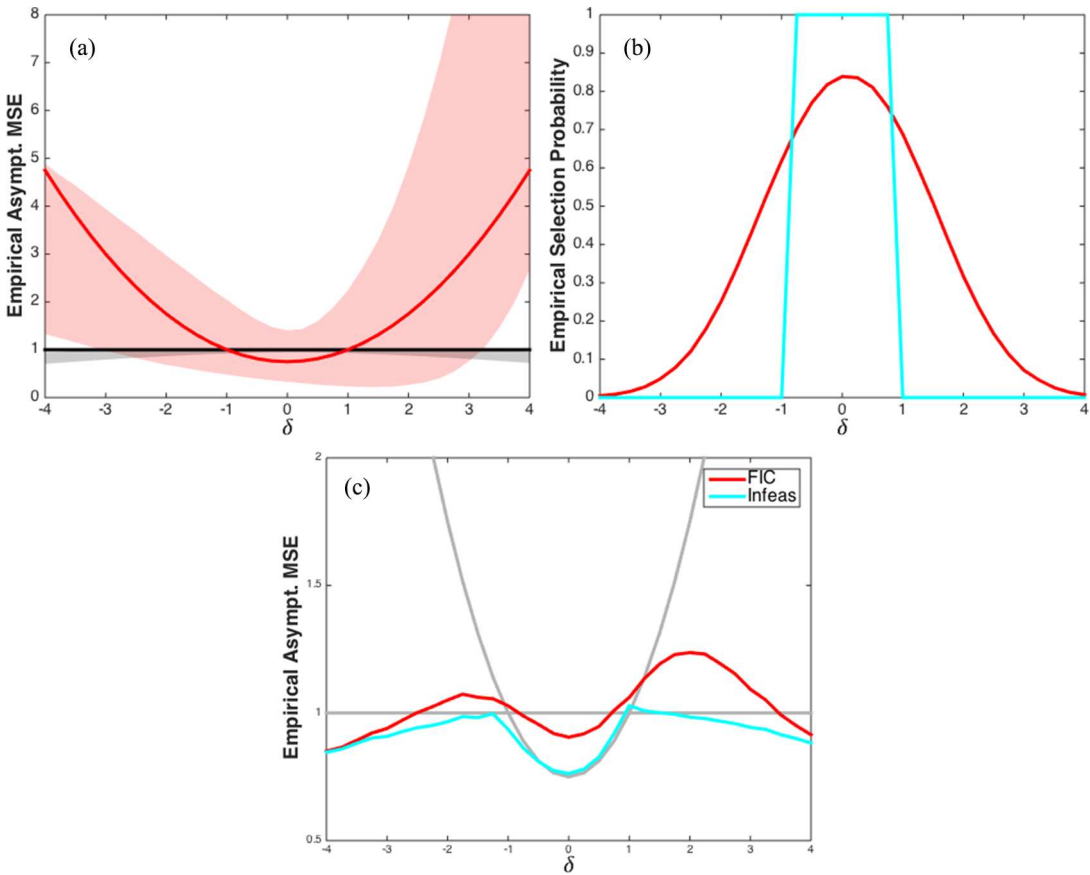
**Figure 3.** (a) The empirical asymptotic MSE of the models with one and two lags (red and black line, respectively). The area between the 5 and 95% empirical quantiles of $\widehat{FIC}_1$ and $\widehat{FIC}_2$ are shaded in red and grey. (b) The empirical selection probabilities of the FIC. (c) The AMSE of the models with $m = 1$ and $m = 2$ together with the empirical MSE of the feasible FIC (red) and infeasible FIC (cyan). This figure was obtained for $T = 100$ and should be compared with the asymptotic results in Figure 1.

### 3.1. Simplified DGP

Figures 3 and 4 provide the finite sample confirmation of the intuition we gained from the simplified DGP.[16] The wide spread in the empirical distribution of $\widehat{FIC}_1$ shown in Figure 3(a) results indeed in a smoothed instead of binary selection between the models (Figure 3(b)). The performance of the feasible FIC is therefore worse than that of the infeasible FIC that assumes $\delta$ to be known. At high $|\delta|$, we see that the probability to select the wrong model is small. The feasible FIC, therefore, performs similarly to its infeasible counterpart for large amounts of misspecification only.

In Figure 4, the three panels display results on the plug-in averages. The quantiles of the weight distribution without bias correction should be compared to those in Equation (2.20) and Figure 2. The results match. Figure 4(b) shows the quantiles of the weight distribution with bias correction. As expected, this distribution is shifted toward higher weights because the upward bias of the AMSE of the model with $m = 1$ is removed. We can see in 4(c) that the plug-in averages do not perform as well as the infeasible estimator. Unreported simulation results at a sample size of $T = 1000$ confirm that this effect does not disappear with sample size. The inconsistent estimation of $\delta$ again causes the feasible weights to differ from optimal weights.

---

[16]In this section we have rescaled the empirical MSE by the sample size to make it comparable to the asymptotic results of Figures 1 and 2, hence the label empirical asymptotic MSE.

**Figure 4.** The 5 and 95% empirical quantiles of the weights distribution without bias correction (a) and with bias correction (b). The infeasible weights are displayed in cyan. The empirical MSE of plug-in methods is shown in (c). The sample size is $T = 100$.

## 3.2. Simulation results for an autoregressive model

Further simulations are based on the following model[17]:

$$y_{T,t} = 0.5y_{T,t-1} + \frac{\delta}{\sqrt{T}}y_{T,t-2} + \frac{\delta}{2\sqrt{T}}y_{T,t-3} + u_t, \qquad u_t \overset{i.i.d.}{\sim} \mathcal{N}(0,1). \tag{3.1}$$

The second and third lag are local-to-zero implying that $\mathcal{M} = \{1, 2, 3\}$. The coefficients in front of the misspecified lags decline linearly as in Liu (2015), where $\delta$ governs the amount of misspecification. The largest modulus eigenvalue of the companion matrix is about 0.3 at $\delta = -0.2$ and increases monotonously to approximately 0.9 at the boundaries of the interval $[-4, 2]$.

**Remark 10.** The AMSEs of the impulse response at horizon 1 are the same for $m = 2$ and $m = 3$. The plug-in weights are not unique, also see Remark 6.

### 3.2.1. MSE comparison
The empirical MSE of the various selection methods are shown in Figure 5 for the impulse responses at horizon 1, 3, and 5. Due to the strong penalty on model complexity, the BIC performs well for small

---

[17]We show in the Supplementary Material that the simplified model with $p_1 = p_2 = 1$, i.e., $y_{T,t} = \alpha y_{T,t-1} + \frac{\delta}{\sqrt{T}}y_{T,t-2} + u_t$, is

special because the gradient vector has no influence on model selection and plug-in averaging. We extend the model with an additional lag to see the influence of the impulse response horizon.

**Figure 5.** The empirical MSE for model selection. The DGP is $y_{T,t} = 0.5y_{T,t-1} + \frac{\delta}{\sqrt{T}}y_{T,t-2} + \frac{\delta}{2\sqrt{T}}y_{T,t-3} + u_t$. (a) $h = 1, T = 100$, (b) $h = 1, T = 1000$, (c) $h = 3, T = 100$, (d) $h = 3, T = 1000$, (e) $h = 5, T = 100$, and (f) $h = 5, T = 1000$.

amounts of misspecification, but its performance quickly deteriorates as $|\delta|$ increases. The performance of the AIC and the feasible version of the FIC are comparable for large areas of the parameter space, with neither of these methods being preferred to the other. The infeasible FIC is very frequently the preferred method.

Model averaging results are reported in Figure 6. The behavior of the smoothed BIC procedure is similar to that of its selection counterpart, i.e., it only performs well for small $\delta$. The same remark applies to the plug-in average with bias correction. The Jackknife, smoothed AIC, and the plug-in average without bias correction are close competitors, where the plug-in average is a better candidate for



**Figure 6.** The empirical MSE for model averaging. The DGP is $y_{T,t} = 0.5y_{T,t-1} + \frac{\delta}{\sqrt{T}}y_{T,t-2} + \frac{\delta}{2\sqrt{T}}y_{T,t-3} + u_t$. (a) $h = 1, T = 100$, (b) $h = 1, T = 1000$, (c) $h = 3, T = 100$, (d) $h = 3, T = 1000$, (e) $h = 5, T = 100$, and (f) $h = 5, T = 1000$.

**Figure 7.** (a) and (b) The empirical MSE of the OLS estimator of the model with 1 lag (OLS1), 2 lags (OLS2) and the full model with 3 lags (OLS3). Gray lines show the asymptotic MSE approximations as provided by the delta method. (c) and (d) The empirical selection probabilities (see Figure 5 for the appropriate legend). (e) and (f) The empirical distribution of the weights (see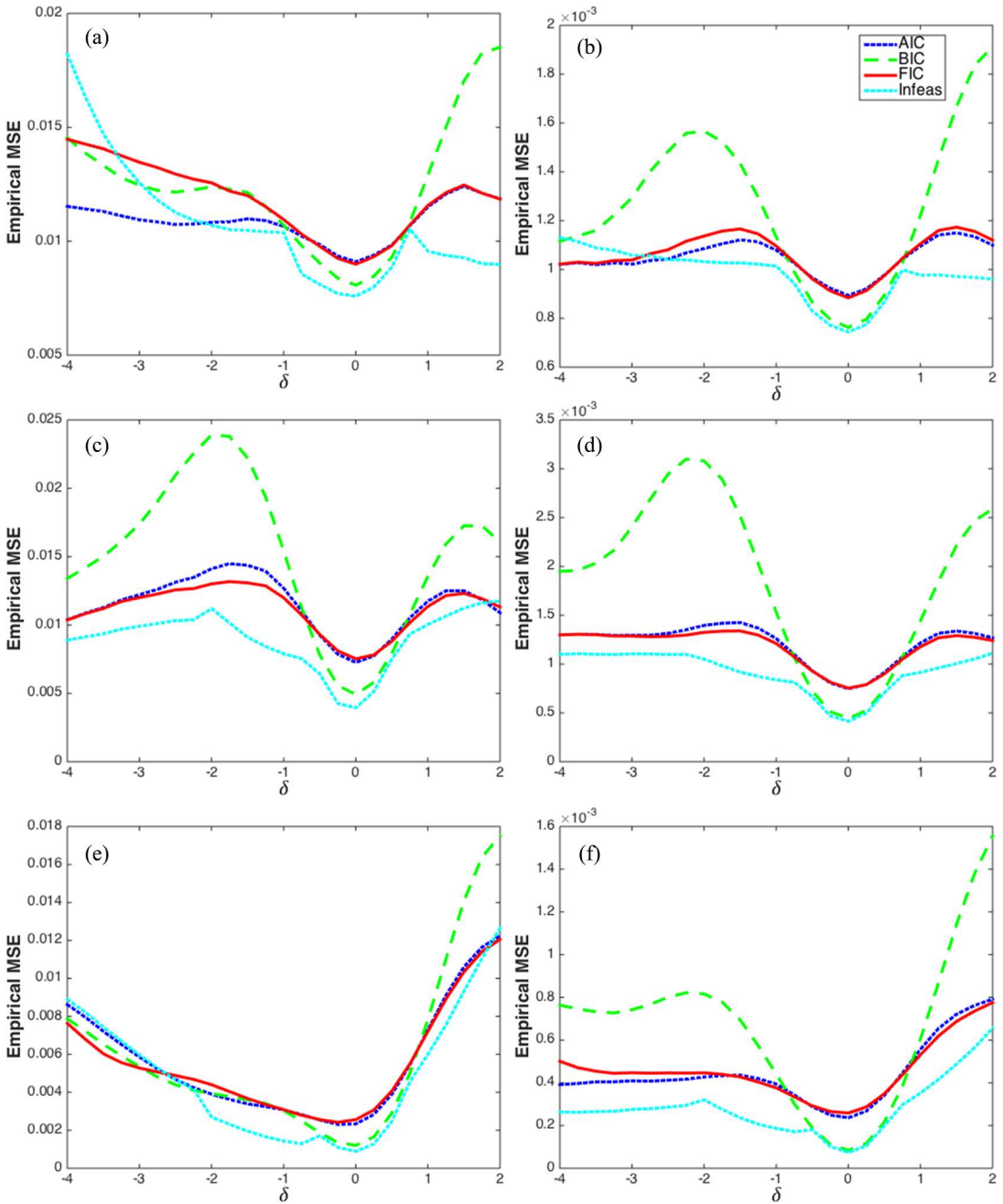 Figure 6 for the appropriate legend). The DGP is $y_{T,t} = 0.5y_{T,t-1} + \frac{\delta}{\sqrt{T}}y_{T,t-2} + \frac{\delta}{2\sqrt{T}}y_{T,t-3} + u_t$ for all graphs. (a) $h = 3$, $T = 100$, (b) $h = 3$, $T = 1000$, (c) $h = 3$, $T = 100$, (d) $h = 3$, $T = 1000$, (e) $h = 3$, $T = 100$, and (f) $h = 3$, $T = 1000$.

Table 1. The empirical coverage of 90% confidence intervals for the horizons: two, three and six. For several choices of $\delta$ and $T$ in the DGP $y_{T,t} = 0.5 y_{T,t-1} + \frac{\delta}{\sqrt{T}} y_{T,t-2} + \frac{\delta}{2\sqrt{T}} y_{T,t-3} + u_t$.

| | h = 2 | | | | h = 3 | | | | h = 6 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\delta$ | T = 100 | T = 250 | T = 500 | T = 1000 | T = 100 | T = 250 | T = 500 | T = 1000 | T = 100 | T = 250 | T = 500 | T = 1000 |
| −4.0 | 75.76 | 81.70 | 83.04 | 84.08 | 76.54 | 81.13 | 81.99 | 82.31 | 75.85 | 58.81 | 42.44 | 2.04 |
| −3.0 | 80.06 | 82.92 | 83.72 | 84.40 | 79.82 | 82.08 | 82.42 | 82.42 | 67.85 | 47.84 | 6.19 | 11.01 |
| −2.0 | 82.33 | 83.79 | 84.16 | 84.76 | 81.92 | 82.64 | 82.54 | 82.46 | 55.30 | 12.54 | 14.52 | 23.33 |
| −1.0 | 83.59 | 84.44 | 84.50 | 84.99 | 82.23 | 82.38 | 82.13 | 82.07 | 24.76 | 23.12 | 27.72 | 31.39 |
| 0.0 | 83.92 | 84.59 | 84.61 | 85.00 | 81.01 | 81.52 | 81.39 | 81.43 | 53.66 | 46.20 | 44.24 | 43.45 |
| 1.0 | 83.20 | 84.33 | 84.51 | 84.94 | 78.11 | 79.67 | 80.13 | 80.58 | 68.96 | 62.81 | 57.59 | 53.38 |
| 2.0 | 80.60 | 83.49 | 84.20 | 84.82 | 71.70 | 76.94 | 78.61 | 79.59 | 76.69 | 73.19 | 67.54 | 61.70 |

large $|\delta|$. The performance of the plug-in average with known $\delta$ is best. It even performs uniformly the best at the larger sample size of $T = 1000$.

What causes the superior performance of the infeasible estimators? Our simulation findings can be understood from the intuition that was gained from the simplified DGP. Panel (a) and (b) from Figure 7 show the empirical MSE of the three models, $m \in \{1, 2, 3\}$, together with the AMSE of these models. The asymptotic approximation is close for $T = 100$ and improves further at $T = 1000$. The selection probabilities in panels (c) and (d) reveal how the infeasible estimator takes a binary decision with the lag length increasing with $\delta$. For the simplified DGP we have seen how the convergence in distribution of $\hat{\delta}$ causes smeared out selection probabilities instead of the binary decision. This effect is also observed in the graphs, *even at the large sample size of* $T = 1000$. The panels (e) and (f) tell the same story for the plug-in weights.

We also perform simulation where we focus on several impulse responses simultaneously, see Remarks 4 and 5. The trace is used to map the AMSE matrix to a scalar. The simulation outcomes are qualitatively similar to our results for the impulse responses at a single horizon. Further details can be found in the Supplementary Material.

### 3.2.2. Confidence intervals

Confidence intervals/bands can be calculated based on Theorem 4. Simulation results are provided in Tables 1 and 2. The desired nominal coverage level was 90%. Table 1 shows that the empirical coverage of the individual confidence intervals is consistently too low. At horizon 2 and 3 this under-coverage is least severe and decreases with sample size. The coverage of the confidence level for horizon 6 varies strongly across $\delta$ and can be very low. It is well-established in the literature (e.g., Kilian (1998) and Kilian

**Table 2.** The empirical coverage of the 90% confidence regions (see Theorem 4) as a function of sample size $T$ and misspecification parameter $\delta$.

| | $h = \{2,3\}$ | | | | $h = \{2,6\}$ | | | |
|---|---|---|---|---|---|---|---|---|
| $\delta$ | $T=100$ | $T=250$ | $T=500$ | $T=1000$ | $T=100$ | $T=250$ | $T=500$ | $T=1000$ |
| −4.0 | 89.71 | 89.94 | 90.01 | 90.08 | 87.14 | 86.47 | 71.25 | 9.87 |
| −3.0 | 89.59 | 90.03 | 90.11 | 90.11 | 84.15 | 73.27 | 41.74 | 56.54 |
| −2.0 | 89.64 | 90.05 | 90.07 | 90.02 | 74.91 | 56.85 | 68.57 | 93.01 |
| −1.0 | 89.59 | 90.01 | 90.01 | 90.02 | 76.13 | 89.96 | 90.41 | 90.21 |
| 0.0 | 89.41 | 89.86 | 89.89 | 89.93 | 83.80 | 86.91 | 88.31 | 89.16 |
| 1.0 | 89.17 | 89.74 | 89.77 | 89.87 | 85.99 | 87.91 | 88.56 | 89.13 |
| 2.0 | 88.54 | 89.75 | 89.79 | 89.91 | 87.29 | 88.83 | 89.16 | 89.46 |



**Figure 8.** The confidence intervals are based on the asymptotic normality of $\sqrt{T}\left(\bar{\mu}(\hat{w}) - \mu(\theta_T, \sigma)\right) - D_\theta \sum_{m=p_1}^{p_1+p_2} \hat{w}_m \hat{C}_m \hat{\delta}$ (see Theorem 4). The displayed histograms are constructed for $y_{T,t} = 0.5y_{T,t-1} + \frac{\delta}{\sqrt{T}}y_{T,t-2} + \frac{\delta}{2\sqrt{T}}y_{T,t-3} + u_t$ with $\delta = -2$ and $T = 500$, i.e., the boxed entries in Table 1. The number of Monte Carlo replications is 100,000. (a) $h = 2, T = 500$ and (b) $h = 6, T = 500$.

(2001)) that inference on impulse responses at higher horizons is inherently more difficult because of the increased nonlinearity in the parameters. This nonlinearity causes the delta method approximation to perform poorly. Figure 8 shows the histograms of $\sqrt{T}\left(\bar{\mu}(\hat{w}) - \mu(\theta_T, \sigma)\right) - D_\theta \sum_{m=p_1}^{p_1+p_2} \hat{w}_m \hat{C}_m \hat{\delta}$ for the impulse responses at horizons 2 and 6 (corresponding to the boxed numbers in Table 1). Note that the confidence intervals/bands defined in Theorem 4 are based on the asymptotic normality of this expression. The sometimes severe under-coverage at horizon 6 should therefore not come as a surprise. This poor asymptotic approximation at horizon 6 also influences the empirical coverage of the confidence bands as can be seen in Table 2.

### 3.3. Simulation results for a vector autoregressive model

Our simulation results are based on a bivariate VAR with DGP

$$
\mathbf{y}_{T,t} = \begin{pmatrix} 0.5 & 0 \\ 0.5 & 0.5 \end{pmatrix} \mathbf{y}_{T,t-1} + \frac{\delta}{\sqrt{T}} \begin{pmatrix} 1 & 0 \\ 0.5 & 1 \end{pmatrix} \mathbf{y}_{T,t-2} + \frac{\delta}{2\sqrt{T}} \begin{pmatrix} 1 & 0 \\ 0.5 & 1 \end{pmatrix} \mathbf{y}_{T,t-3} + \mathbf{u}_t,
$$

$$
\mathbf{u}_t \overset{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_u), \text{ where } \boldsymbol{\Sigma}_u = \begin{pmatrix} 1 & 0.17 \\ 0.17 & 0.33 \end{pmatrix},
$$

(3.2)



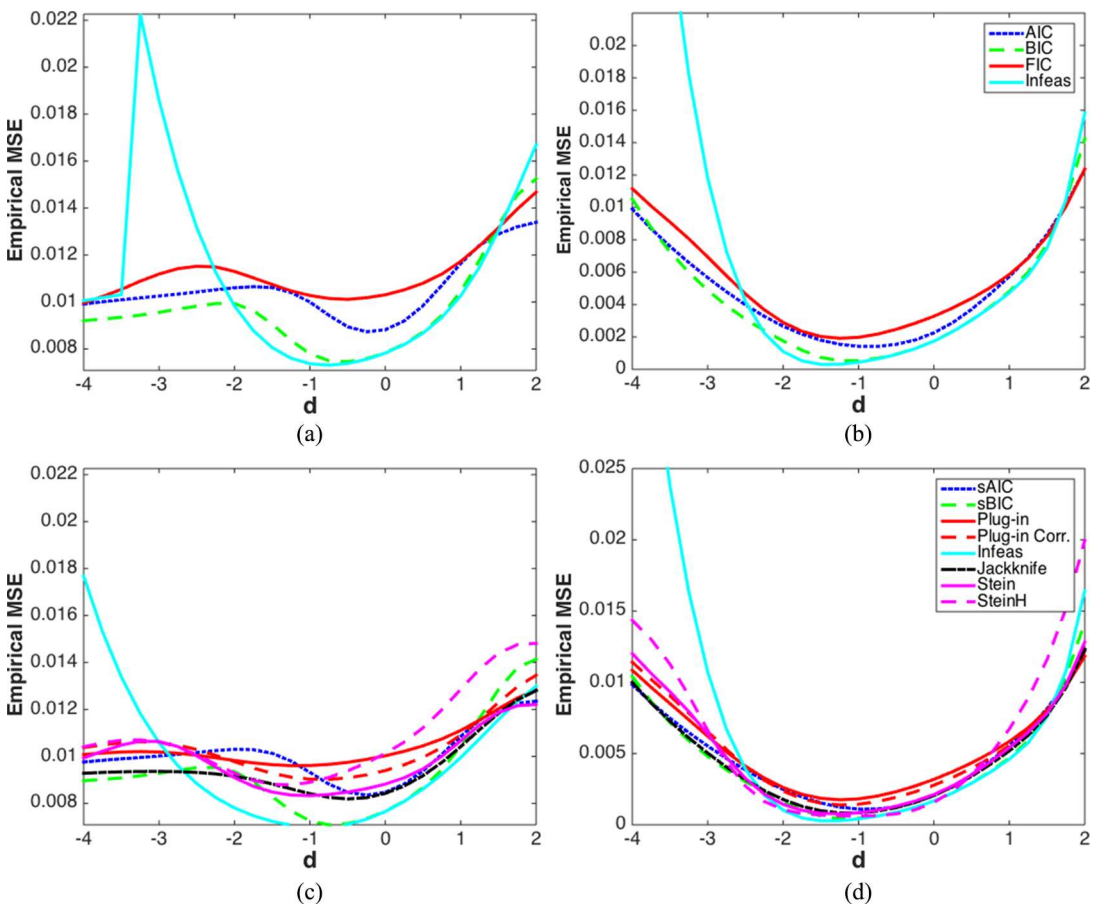**Figure 9.** The empirical MSE of the impulse response estimator for several selection and averaging methods. We have displayed the results for the response of variable 1 to a structural shock in the variable 1 for horizons 2 and 6. The DGP is given in Equation (3.2). The sample size is $T = 100$. (a) $h = 2$, model selection, (b) $h = 6$, model selection, (c) $h = 2$, model averaging, and (d) $h = 6$, model averaging.

**Figure 10.** Identical to Figure 9, but for $T = 1000$. (a) $h = 2$, model selection, (b) $h = 6$, model selection, (c) $h = 2$, model averaging, and (d) $h = 6$, model averaging.

which is similar to the VAR used in Lütkepohl et al. (2015) for impulse response analyses. This process has the same roots as the univariate process of Equation (3.1) but with double multiplicity. The parameter $\delta$ governs the degree of misspecification. For brevity, we only report MSE results of the response of variable 1 to a structural shock in variable 1. Figures 9, and 10 show the results for horizons 2 and 6.[18] Similarly to the univariate results, none of the methods performs uniformly best. Only the infeasible methods get close to dominating all other methods for the large sample size of $T = 1000$. The ragged spike for "**Infeas**" in Figure 9(a) is caused by an abrupt binary decision to switch between models with different lag lengths. Finally, it is interesting to note that the Stein shrinkage methods perform well in comparison to the plug-in averaging procedure.

## 4. Conclusion

In this paper, we studied the issue of model selection and model averaging for multivariate autoregressive processes in a locally drifting asymptotic framework. Within this drifting framework we derived the

---

[18]The simulation results for all four impulse responses and horizons 1–6 can be found in the Supplementary Material.

asymptotic normality of the least squares estimators. The multivariate delta method subsequently ensured that this asymptotic normality carries over to sufficiently smooth parameter transformations, e.g., impulse responses. The focused information criterion and plug-in averaging estimator were defined as the minimizers of the estimated asymptotic mean squared error of the focus parameter estimator.

We highlighted the role of the misspecification parameter $\delta$. Both Liu (2015) and DiTraglia (2016) mentioned that the feasible FIC remains random in the limit. We provided the explicit expressions for the limiting distribution of the FIC values and the elements of the weighting matrices, and illustrated that the feasible estimators do not truly minimize the asymptotic mean squared error. This latter result might encourage further research into different ways to deal with the misspecification parameter. There are to the best of our knowledge two alternatives reported in the literature. The recent paper by Kitagawa and Muris (2016) adopts a mixed frequentist and Bayesian framework to alleviate the estimation of $\delta$ in their study of model averaging in semiparametric estimation of treatment effects. Hansen (2016) similarly adopts a local-to-zero framework but minimizes a risk quantity that does not require the direct estimation of $\delta$.

Our simulation study of univariate and multivariate autoregressive processes underlined the previous paragraph because the infeasible estimator (the estimator knowing $\delta$) frequently dominated the other methods. The latter was especially the case at the larger sample size of $T = 1000$. There was no clearly preferred method for *feasible* model selection/averaging.

A possible extension of this work is an application to forecasting. Such an extension would complement: (1) the predictive static regression setup discussed in Liu and Kuo (2016), and (2) the prediction focused model selection of autoregressive models in Claeskens et al. (2007). Forecasts for autoregressive models often start from the assumption that estimation and prediction are applied to two independent processes with the same stochastic structure. The link to this current paper is that (under this independence assumption) the asymptotic covariance matrix of the forecast is a continuous transformation of the autoregressive parameters, see Section 3.5 of Lütkepohl (2005) for further details.

## Appendix

## Mathematical proofs

*Proof of Theorem 1.* As in Liu (2015) we first relate the parameters of the VAR($m$) models to the parameters of the VAR($p$). With the aid of the selection matrices we have

$$
\begin{aligned}
\hat{\boldsymbol{\Theta}}_{T,m} &= \boldsymbol{Y}_T \boldsymbol{Z}'_{T,m} (\boldsymbol{Z}_{T,m} \boldsymbol{Z}'_{T,m})^{-1} = (\boldsymbol{BL}' \boldsymbol{Z}_T + \boldsymbol{C}_T \boldsymbol{S}'_0 \boldsymbol{Z}_T + \boldsymbol{U}) \boldsymbol{Z}'_{T,m} (\boldsymbol{Z}_{T,m} \boldsymbol{Z}'_{T,m})^{-1} \\
&= ([\boldsymbol{B} \; \boldsymbol{C}_T] \boldsymbol{S}_m \boldsymbol{S}'_m \boldsymbol{Z}_T + \boldsymbol{C}_T (\boldsymbol{I}_{Kp_2} - \boldsymbol{\Pi}'_m \boldsymbol{\Pi}_m) \boldsymbol{S}'_0 \boldsymbol{Z}_T + \boldsymbol{U}) \boldsymbol{Z}'_{T,m} (\boldsymbol{Z}_{T,m} \boldsymbol{Z}'_{T,m})^{-1} \\
&= (\boldsymbol{\Theta}_{T,m} \boldsymbol{Z}_{T,m} + \boldsymbol{C}_T (\boldsymbol{I}_{Kp_2} - \boldsymbol{\Pi}'_m \boldsymbol{\Pi}_m) \boldsymbol{S}'_0 \boldsymbol{Z}_T + \boldsymbol{U}) \boldsymbol{Z}'_{T,m} (\boldsymbol{Z}_{T,m} \boldsymbol{Z}'_{T,m})^{-1} \\
&= \boldsymbol{\Theta}_{T,m} + \boldsymbol{C}_T (\boldsymbol{I}_{Kp_2} - \boldsymbol{\Pi}'_m \boldsymbol{\Pi}_m) \boldsymbol{S}'_0 \boldsymbol{Z}_T \boldsymbol{Z}'_{T,m} (\boldsymbol{Z}_{T,m} \boldsymbol{Z}'_{T,m})^{-1} + \boldsymbol{U} \boldsymbol{Z}'_{T,m} (\boldsymbol{Z}_{T,m} \boldsymbol{Z}'_{T,m})^{-1}. \quad \text{(A.1)}
\end{aligned}
$$

We rearrange terms to obtain the starting point of our analysis. Note especially how the scaling by $T^{1/2}$ cancels against the $T^{-1/2}$ decay rate of the elements in the matrix $C_T$. See Remark 1 for a discussion. An expression in terms of the fixed parameter matrix $\boldsymbol{\Delta}$ remains:

$$
\sqrt{T}\left(\hat{\boldsymbol{\Theta}}_{T,m} - \boldsymbol{\Theta}_{T,m}\right) = \underbrace{\sqrt{T} \boldsymbol{C}_T \left(\boldsymbol{I}_{Kp_2} - \boldsymbol{\Pi}'_m \boldsymbol{\Pi}_m\right) \boldsymbol{S}'_0}_{\boldsymbol{\Delta}} \left(\frac{1}{T} \boldsymbol{Z}_T \boldsymbol{Z}'_T\right) \boldsymbol{S}_m \left[\boldsymbol{S}'_m \left(\frac{1}{T} \boldsymbol{Z}_T \boldsymbol{Z}'_T\right) \boldsymbol{S}_m\right]^{-1}
$$

$$
+ \left(\frac{1}{\sqrt{T}} \boldsymbol{U} \boldsymbol{Z}'_T\right) \boldsymbol{S}_m \left[\boldsymbol{S}'_m \left(\frac{1}{T} \boldsymbol{Z}_T \boldsymbol{Z}'_T\right) \boldsymbol{S}_m\right]^{-1}. \quad \text{(A.2)}
$$

If we define the vector of parameter estimates as $\hat{\boldsymbol{\theta}}_{T,m} = \text{vec}(\hat{\boldsymbol{\Theta}}_{T,m})$ and the true parameter vector $\boldsymbol{\theta}_{T,m} = \text{vec}(\boldsymbol{\Theta}_{T,m})$, then the properties of the vec operator provide

$$\sqrt{T}\left(\hat{\boldsymbol{\theta}}_{T,m} - \boldsymbol{\theta}_{T,m}\right) = \left(\left[\boldsymbol{S}_m'\left(\frac{1}{T}\boldsymbol{Z}_T\boldsymbol{Z}_T'\right)\boldsymbol{S}_m\right]^{-1}\boldsymbol{S}_m'\left(\frac{1}{T}\boldsymbol{Z}_T\boldsymbol{Z}_T'\right)\boldsymbol{S}_0\left(\boldsymbol{I}_{Kp_2} - \boldsymbol{\Pi}_m'\boldsymbol{\Pi}_m\right)\otimes\boldsymbol{I}_K\right)\boldsymbol{\delta}$$

$$+\left(\left[\boldsymbol{S}_m'\left(\frac{1}{T}\boldsymbol{Z}_T\boldsymbol{Z}_T'\right)\boldsymbol{S}_m\right]^{-1}\boldsymbol{S}_m'\otimes\boldsymbol{I}_K\right)\frac{1}{\sqrt{T}}\sum_{t=1}^{T}\text{vec}\left(\boldsymbol{u}_t\boldsymbol{z}_{T,t-1}'\right), \qquad \text{(A.3)}$$

where $\boldsymbol{\delta} = \text{vec}(\boldsymbol{\Delta})$. We will prove both $\text{plim}_{T\to\infty}\frac{1}{T}\boldsymbol{Z}_T\boldsymbol{Z}_T' = \text{plim}_{T\to\infty}\frac{1}{T}\sum_{t=1}^{T}\boldsymbol{z}_{T,t-1}\boldsymbol{z}_{T,t-1}' = \boldsymbol{\Omega}$ and $\frac{1}{\sqrt{T}}\sum_{t=1}^{T}\text{vec}(\boldsymbol{u}_t\boldsymbol{z}_{T,t-1}') \xrightarrow{d} \text{N}(\boldsymbol{0},\boldsymbol{\Omega}\otimes\boldsymbol{\Sigma})$. These two results prove part (a) of Theorem 1, because the continuous mapping theorem implies that $\sqrt{T}(\hat{\boldsymbol{\theta}}_{T,m} - \boldsymbol{\theta}_{T,m}) = \boldsymbol{A}_m\boldsymbol{\delta} + ([\boldsymbol{S}_m'\boldsymbol{\Omega}\boldsymbol{S}_m]^{-1}\boldsymbol{S}_m'\otimes\boldsymbol{I}_K)\boldsymbol{R} + o_p(1)$.

We start with the proof of $\text{plim}_{T\to\infty}\frac{1}{T}\boldsymbol{Z}_T\boldsymbol{Z}_T' = \boldsymbol{\Omega}$. The process $\{\boldsymbol{y}_{T,t}\}_{t=-\infty}^{\infty}$ is stationary and ergodic for every fixed $T$ in view of Assumptions 1 and 2 (e.g., Theorem 3 on p. 204 of Hannan (1970)). Define the companion matrix $\boldsymbol{A}_T$ and innovation vector $\boldsymbol{E}_t$ such that $\boldsymbol{z}_{T,t} = \boldsymbol{A}_T\boldsymbol{z}_{T,t-1} + \boldsymbol{E}_t$, i.e.,

$$\boldsymbol{A}_T := \begin{bmatrix} \boldsymbol{B}_1 & \boldsymbol{B}_2 & \cdots & \boldsymbol{B}_{p_1-1} & \boldsymbol{B}_{p_1} & \frac{\boldsymbol{C}_1}{\sqrt{T}} & \frac{\boldsymbol{C}_2}{\sqrt{T}} & \cdots & \frac{\boldsymbol{C}_{p_2-1}}{\sqrt{T}} & \frac{\boldsymbol{C}_{p_2}}{\sqrt{T}} \\ \boldsymbol{I}_K & \boldsymbol{O} & \cdots & \boldsymbol{O} & \boldsymbol{O} & \boldsymbol{O} & \boldsymbol{O} & \cdots & \boldsymbol{O} & \boldsymbol{O} \\ \vdots & \vdots & & \vdots & \vdots & \vdots & \vdots & & \vdots & \vdots \\ \boldsymbol{O} & \boldsymbol{O} & \cdots & \boldsymbol{O} & \boldsymbol{O} & \boldsymbol{O} & \boldsymbol{O} & \cdots & \boldsymbol{I}_K & \boldsymbol{O} \end{bmatrix} \quad (Kp \times Kp),$$

$$\boldsymbol{E}_t := \left(\boldsymbol{u}_t', \boldsymbol{0}', \ldots, \boldsymbol{0}'\right)' \qquad (Kp \times 1). \qquad \text{(A.4)}$$

From this extended VAR(1) form we conclude that $\boldsymbol{z}_{T,t}\boldsymbol{z}_{T,t}' = \boldsymbol{A}_T\boldsymbol{z}_{T,t-1}\boldsymbol{z}_{T,t-1}'\boldsymbol{A}_T' + \boldsymbol{A}\boldsymbol{z}_{T,t-1}\boldsymbol{E}_t' + \boldsymbol{E}_{t-1}\boldsymbol{z}_{T,t-1}'\boldsymbol{A}_T' + \boldsymbol{E}_t\boldsymbol{E}_t'$. Stationarity implies

$$\lim_{T\to\infty}\left(\boldsymbol{I}_{K^2p^2} - \boldsymbol{A}_T\otimes\boldsymbol{A}_T\right)\text{plim vec}\left(\frac{1}{T}\sum_{t=1}^{T}\boldsymbol{z}_{T,t-1}\boldsymbol{z}_{T,t-1}'\right)$$

$$= \text{plim vec}\left(\frac{1}{T}\sum_{t=1}^{T}\boldsymbol{A}_T\boldsymbol{z}_{T,t-1}\boldsymbol{E}_t'\right) + \text{plim vec}\left(\frac{1}{T}\sum_{t=1}^{T}\boldsymbol{E}_t\boldsymbol{z}_{T,t-1}'\boldsymbol{A}_T'\right) + \text{plim vec}\left(\frac{1}{T}\sum_{t=1}^{T}\boldsymbol{E}_t\boldsymbol{E}_t'\right). \qquad \text{(A.5)}$$

The nonrandom matrix $\boldsymbol{A}_T$ will converge to the matrix $\boldsymbol{A}_\infty$ for large $T$. $\boldsymbol{A}_\infty$ is thus obtained from $\boldsymbol{A}_T$ by replacing the ratios $\boldsymbol{C}_i/\sqrt{T}$ with zero matrices for $i \in \{1, 2, \ldots, p_2\}$. Note that the eigenvalues of $\boldsymbol{A}_\infty$ coincide with the roots of the matrix polynomial $\boldsymbol{B}_\infty(z)$ augmented with $Kp_2$ additional zero eigenvalues. Assumption 3 guarantees that the matrix $\boldsymbol{I}_{p^2} - \boldsymbol{A}_\infty\otimes\boldsymbol{A}_\infty$ is invertible.

Subsequently we consider the RHS of Equation (A.5). Let $y_{T,t-j,k}$ and $u_{t,k}$ denote the $k$th component of $\boldsymbol{y}_{T,t-j}$ and $\boldsymbol{u}_t$, respectively. If we can show that $\frac{1}{T}\sum_{t=1}^{T}y_{T,t-j,k}u_{t,l} \xrightarrow{p} 0$ for all $j \in \{1, 2, \ldots, p\}$ and $k, l \in \{1, 2, \ldots, K\}$, then the first two terms in the RHS of Equation (A.5) are $o_p(1)$. To prove this, we define the array $X_{T,t}^{jkl} = y_{T,t-j,k}u_{t,l}/T$ and the norming $c_T = 1/T$. $X_{T,t}^{jkl}$ is a martingale difference (m.d.) array with respect to the filtration $\mathcal{F}_t = \sigma(\boldsymbol{u}_s, -\infty < s \le t)$ and $\text{E}|X_{T,t}^{jkl}/c_T|^4 = \text{E}|y_{T,t-j,k}u_{t,l}|^4$ is finite in view of Assumption 1. Result 12.10 from Davidson (1994) implies that $|X_{T,t}^{jkl}|^2$ is uniformly integrable and Result 19.7 from the same reference gives $\frac{1}{T}\sum_{t=1}^{T}y_{T,t-j,k}u_t \xrightarrow{L_2} 0$. The result for the first two terms follows. The third term in the RHS of Equation (A.5) is a sample mean of an i.i.d. sequence. Khinchine's

Theorem gives the probability limit. Combining all the results, we conclude that

$$\text{vec}\,(\boldsymbol{\Omega}) := \text{plim}_{T\to\infty}\text{vec}\left(\frac{1}{T}\sum_{t=1}^{T}\boldsymbol{z}_{T,t-1}\boldsymbol{z}'_{T,t-1}\right) = \left(\boldsymbol{I}_{K^2p^2} - \boldsymbol{A}_\infty \otimes \boldsymbol{A}_\infty\right)^{-1}\text{vec}\left(\boldsymbol{\Sigma}^*\right) + o_p(1),\quad \text{(A.6)}$$

where $\boldsymbol{\Sigma}^* = \boldsymbol{e}\boldsymbol{\Sigma}\boldsymbol{e}'$ and $\boldsymbol{e}$ is the $(Kp \times K)$ matrix composed of the first $K$ column of $I_{Kp}$. This shows that $\text{plim}_{T\to\infty}\frac{1}{T}\boldsymbol{Z}_T\boldsymbol{Z}'_T = \boldsymbol{\Omega}$. Ergodicity for every $T$ also provides the result $(\boldsymbol{I}_{K^2p^2} - \boldsymbol{A}_T \otimes \boldsymbol{A}_T)\text{vec}(\text{E}(\boldsymbol{z}_{T,t-1}\boldsymbol{z}'_{T,t-1})) = \text{vec}\,(\boldsymbol{\Sigma})$ and hence $\boldsymbol{\Omega} = \lim_{T\to\infty}\text{E}(\boldsymbol{z}_{T,t-1}\boldsymbol{z}'_{T,t-1})$ because $\boldsymbol{A}_T \to \boldsymbol{A}_\infty$.

We rely on the Cramer-Wold theorem (e.g., Result 25.5 from Davidson (1994)) to prove the convergence of $\frac{1}{\sqrt{T}}\sum_{t=1}^{T}\text{vec}(\boldsymbol{u}_t\boldsymbol{z}'_{T,t-1})$ to $\boldsymbol{R} \sim \text{N}\,(\boldsymbol{0}, \boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})$. Let $\boldsymbol{\xi}$ denote a fixed $(K^2p \times 1)$ vector. $X^*_{T,t} = \boldsymbol{\xi}'\text{vec}(\boldsymbol{u}_t\boldsymbol{z}'_{T,t-1})$ is a m.d. array with respect to $\mathcal{F}_t$. We note that $\sigma^2_{Tt} = \text{E}(X^{*2}_{T,t}|\mathcal{F}_{t-1}) = \text{E}(\boldsymbol{\xi}'\text{vec}(\boldsymbol{u}_T\boldsymbol{z}'_{T,t-1})\text{vec}(\boldsymbol{u}_T\boldsymbol{z}'_{T,t-1})'\boldsymbol{\xi}|\mathcal{F}_{t-1}) = \text{E}(\boldsymbol{\xi}'((\boldsymbol{z}_{T,t-1}\boldsymbol{z}'_{T,t-1})\otimes(\boldsymbol{u}_t\boldsymbol{u}_t))\boldsymbol{\xi}|\mathcal{F}_{t-1}) = \boldsymbol{\xi}'((\boldsymbol{z}_{T,t-1}\boldsymbol{z}'_{T,t-1})\otimes \boldsymbol{\Sigma})\boldsymbol{\xi}$ and that $X^{*2}_{T,t}$ is square integrable by Assumption 1. Moreover, from $s^2_T = \sum_{t=1}^{T}\text{E}(X^{*2}_{T,t}) = \boldsymbol{\xi}'(\text{E}(\boldsymbol{z}_{T,t-1}\boldsymbol{z}'_{T,t-1}) \otimes \boldsymbol{\Sigma})\boldsymbol{\xi}$ we have

$$\sup_T \frac{T}{s^2_T} = \sup_T \frac{1}{\boldsymbol{\xi}'\left(\text{E}\left(\boldsymbol{z}_{T,t-1}\boldsymbol{z}'_{T,t-1}\right)\otimes\boldsymbol{\Sigma}\right)\boldsymbol{\xi}} < \infty. \quad \text{(A.7)}$$

Equation (A.7) holds because the quadratic form in the denominator cannot be zero as both $\boldsymbol{\Sigma}$ and $\text{E}(\boldsymbol{z}_{T,t-1}\boldsymbol{z}'_{T,t-1})$ are positive definite matrices. The positive definiteness of $\boldsymbol{\Sigma}$ is part of Assumption 4. For finite $T$, $\text{E}(\boldsymbol{z}_{T,t-1}\boldsymbol{z}'_{T,t-1})$ cannot be positive semidefinite because this would imply the existence of a $\boldsymbol{\kappa} \neq \boldsymbol{0}$ such that $\boldsymbol{\kappa}'\text{E}(\boldsymbol{z}_{T,t-1}\boldsymbol{z}'_{T,t-1})\boldsymbol{\kappa} = \text{E}(\boldsymbol{\kappa}'\boldsymbol{z}_{T,t-1})^2 = 0$ and at least one component of $\boldsymbol{z}_{T,t-1}$ is zero for all $t$. Also it cannot approach a positive semidefinite matrix due to convergence to $\boldsymbol{\Omega}$. A generalization of Result 24.4 from Davidson (1994) to martingale difference arrays shows that

$$\frac{\frac{1}{\sqrt{T}}\sum_{t=1}^{T}\boldsymbol{\xi}'\text{vec}\left(\boldsymbol{u}_t\boldsymbol{z}'_{T,t-1}\right)}{\sqrt{\boldsymbol{\xi}'\left(\text{E}\left(\boldsymbol{z}_{T,t-1}\boldsymbol{z}'_{T,t-1}\right)\otimes\boldsymbol{\Sigma}\right)\boldsymbol{\xi}}} \xrightarrow{d} \text{N}\,(0,1). \quad \text{(A.8)}$$

The expression under the square root is asymptotically equivalent to $\boldsymbol{\xi}'(\boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})\boldsymbol{\xi}$. The second result, $\frac{1}{\sqrt{T}}\sum_{t=1}^{T}\text{vec}(\boldsymbol{u}_t\boldsymbol{z}'_{T,t-1}) \xrightarrow{d} \text{N}\,(\boldsymbol{0}, \boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})$, follows because $\boldsymbol{\xi}$ is arbitrary. The proof of part (a) is complete.

Part (b) of Theorem 1 is a joint convergence result with the estimator for the covariance matrix. For any model with $m \in \{p_1, p_1 + 1, \ldots, p\}$ lags we define the residual matrix (residuals are stacked columnwise) by $\hat{\boldsymbol{U}}^m_{T,t} = \boldsymbol{Y}_T - \hat{\boldsymbol{B}}_{T,m}\boldsymbol{Z}_{T,m}$. The estimated covariance matrix based on the residuals from the model with $m$ lags satisfies

$$\hat{\boldsymbol{\Sigma}}^m_u = \frac{1}{T}\left(\boldsymbol{Y}_T - \hat{\boldsymbol{\Theta}}_{T,m}\boldsymbol{Z}_{T,m}\right)\left(\boldsymbol{Y}_T - \hat{\boldsymbol{\Theta}}_{T,m}\boldsymbol{Z}_{T,m}\right)'$$

$$= \frac{1}{T}\left[\left(\boldsymbol{\Theta}_{T,m} - \hat{\boldsymbol{\Theta}}_{T,m}\right)\boldsymbol{S}'_m\boldsymbol{Z}_T + \boldsymbol{C}_T\left(\boldsymbol{I}_{Kp_2} - \boldsymbol{\Pi}'_m\boldsymbol{\Pi}_m\right)\boldsymbol{S}'_0\boldsymbol{Z}_T + \boldsymbol{U}\right]$$

$$\left[\left(\boldsymbol{\Theta}_{T,m} - \hat{\boldsymbol{\Theta}}_{T,m}\right)\boldsymbol{S}'_m\boldsymbol{Z}_T + \boldsymbol{C}_T\left(\boldsymbol{I}_{Kp_2} - \boldsymbol{\Pi}'_m\boldsymbol{\Pi}_m\right)\boldsymbol{S}'_0\boldsymbol{Z}_T + \boldsymbol{U}\right]'$$

$$= \left(\boldsymbol{\Theta}_{T,m} - \hat{\boldsymbol{\Theta}}_{T,m}\right)\boldsymbol{S}'_m\left(\frac{\boldsymbol{Z}_T\boldsymbol{Z}'_T}{T}\right)\boldsymbol{S}_m\left(\boldsymbol{\Theta}_{T,m} - \hat{\boldsymbol{\Theta}}_{T,m}\right)'$$

$$+ \boldsymbol{C}_T\left(\boldsymbol{I}_{Kp_2} - \boldsymbol{\Pi}'_m\boldsymbol{\Pi}_m\right)\boldsymbol{S}'_0\left(\frac{\boldsymbol{Z}_T\boldsymbol{Z}'_T}{T}\right)\boldsymbol{S}_0\left(\boldsymbol{I}_{Kp_2} - \boldsymbol{\Pi}'_m\boldsymbol{\Pi}_m\right)'\boldsymbol{C}'_T$$

$$+ \left(\boldsymbol{\Theta}_{T,m} - \hat{\boldsymbol{\Theta}}_{T,m}\right) \boldsymbol{S}'_m \left(\frac{\boldsymbol{Z}_T \boldsymbol{Z}'_T}{T}\right) \boldsymbol{S}_0 \left(\boldsymbol{I}_{Kp_2} - \boldsymbol{\Pi}'_m \boldsymbol{\Pi}_m\right)' \boldsymbol{C}'_T + \left(\boldsymbol{\Theta}_{T,m} - \hat{\boldsymbol{\Theta}}_{T,m}\right) \boldsymbol{S}'_m \left(\frac{\boldsymbol{Z}_T \boldsymbol{U}'}{T}\right)$$

$$+ \boldsymbol{C}_T \left(\boldsymbol{I}_{Kp_2} - \boldsymbol{\Pi}'_m \boldsymbol{\Pi}_m\right) \boldsymbol{S}'_0 \left(\frac{\boldsymbol{Z}_T \boldsymbol{Z}'_T}{T}\right) \boldsymbol{S}_m \left(\boldsymbol{\Theta}_{T,m} - \hat{\boldsymbol{\Theta}}_{T,m}\right)' + \boldsymbol{C}_T \left(\boldsymbol{I}_{Kp_2} - \boldsymbol{\Pi}'_m \boldsymbol{\Pi}_m\right) \boldsymbol{S}'_0 \left(\frac{\boldsymbol{Z}_T \boldsymbol{U}'}{T}\right)$$

$$+ \left(\frac{\boldsymbol{U}\boldsymbol{Z}'_T}{T}\right) \boldsymbol{S}_m \left(\boldsymbol{\Theta}_{T,m} - \hat{\boldsymbol{\Theta}}_{T,m}\right) + \left(\frac{\boldsymbol{U}\boldsymbol{Z}_T}{T}\right) \boldsymbol{S}_0 \left(\boldsymbol{I}_{Kp_2} - \boldsymbol{\Pi}'_m \boldsymbol{\Pi}_m\right)' \boldsymbol{C}'_T + \frac{1}{T}\boldsymbol{U}\boldsymbol{U}'. \tag{A.9}$$

The stochastic orders of the various terms in Equation (A.9) are known from previous results. We have $\boldsymbol{\Theta}_{T,m} - \hat{\boldsymbol{\Theta}}_{T,m} = O_p(T^{-1/2})$, $\boldsymbol{Z}_T \boldsymbol{Z}'_T/T \xrightarrow{p} \boldsymbol{\Omega}$, $\boldsymbol{C}_T = O(T^{-1/2})$ and $\boldsymbol{Z}_T \boldsymbol{U}'/T = O_p(T^{-1/2})$ by Equation (A.8). We conclude that $\hat{\boldsymbol{\Sigma}}_u^m = \frac{1}{T}\boldsymbol{U}\boldsymbol{U}' + o_P(1)$. Every covariance estimator (every in the sense of for all $m \in \mathcal{M}$) has therefore the same asymptotic distribution as the covariance estimator based on the true innovations.

Joint asymptotic normality of the parameter estimates and the covariance estimator can be obtained along the lines of the proof of Proposition 11.2 of Hamilton (1994). That is, we define

$$\boldsymbol{\lambda}_t = \text{vech} \begin{bmatrix} u_{1t}^2 - \sigma_{11} & u_{1t}u_{2t} - \sigma_{12} & \dots & u_{1t}u_{Kt} - \sigma_{1K} \\ u_{2t}u_{1t} - \sigma_{21} & u_{2t}^2 - \sigma_{22} & \dots & u_{2t}u_{Kt} - \sigma_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ u_{Kt}u_{1t} - \sigma_{K1} & u_{Kt}u_{2t} - \sigma_{K2} & \dots & u_{Kt}^2 - \sigma_{KK} \end{bmatrix}. \tag{A.10}$$

The sequence $\{\boldsymbol{\lambda}_t\}$ is i.i.d. and thus also a martingale difference sequence. One can apply the Cramer-Wold Theorem to the extended martingale difference vector $(\text{vec}(\boldsymbol{u}_t \boldsymbol{z}'_{T,t-1})', \boldsymbol{\lambda}'_t)'$ to show

$$\begin{bmatrix} (1\sqrt{T}) \sum_{t=1}^T \text{vec}(\boldsymbol{u}_t \boldsymbol{z}'_{T,t-1}) \\ (1\sqrt{T}) \sum_{t=1}^T \boldsymbol{\lambda}_t \end{bmatrix} \xrightarrow{d} \text{N}\left(\begin{bmatrix} \boldsymbol{0} \\ \boldsymbol{0} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Xi}_{11} & \boldsymbol{\Xi}_{12} \\ \boldsymbol{\Xi}_{21} & \boldsymbol{\Xi}_{22} \end{bmatrix}\right). \tag{A.11}$$

We already know that $\boldsymbol{\Xi}_{11} = \boldsymbol{\Omega} \otimes \boldsymbol{\Sigma}$. The elements in the covariance matrix $\boldsymbol{\Xi}_{12}$ take the form $\lim_{T\to\infty} \text{E}(u_{k_1t}y_{T,t-j,k_2}(u_{k_3t}u_{k_4t} - \sigma_{k_3k_4}))$. They are zero because $\lim_{T\to\infty} \text{E}(y_{T,t-j,k_2}) = 0$. Finally, $\boldsymbol{\Xi}_{22} = \text{E}\left(\boldsymbol{\lambda}_t \boldsymbol{\lambda}'_t\right)$. The typical elements are $\text{E}((u_{it}u_{jt} - \sigma_{ij})(u_{lt}u_{mt} - \sigma_{lm}))$.

Define two independent random vectors: $\boldsymbol{R} \sim \text{N}(\boldsymbol{0}, \boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})$ and $\boldsymbol{S} \sim \text{N}(\boldsymbol{0}, \boldsymbol{\Xi}_{22})$. We will proof the claim in part (c) of Theorem 1 for the case of three different models indexed by $m_1, m_2, m_3 \in \mathcal{M}$. The proof is immediate, since

$$\begin{bmatrix} \sqrt{T}\left(\hat{\boldsymbol{\theta}}_{T,m_1} - \boldsymbol{\theta}_{T,m_1}\right) \\ \sqrt{T}\left(\hat{\boldsymbol{\theta}}_{T,m_2} - \boldsymbol{\theta}_{T,m_2}\right) \\ \sqrt{T}\text{vech}\left(\hat{\boldsymbol{\Sigma}}_u^{m_3} - \boldsymbol{\Sigma}_u\right) \end{bmatrix} \xrightarrow{d} \begin{bmatrix} \boldsymbol{A}_{m_1} \\ \boldsymbol{A}_{m_2} \\ \boldsymbol{O} \end{bmatrix} \boldsymbol{\delta} + \begin{bmatrix} \left([\boldsymbol{S}'_{m_1}\boldsymbol{\Omega}\boldsymbol{S}_{m_1}]^{-1}\boldsymbol{S}'_{m_1} \otimes \boldsymbol{I}_K\right) & \boldsymbol{O} \\ \left([\boldsymbol{S}'_{m_2}\boldsymbol{\Omega}\boldsymbol{S}_{m_2}]^{-1}\boldsymbol{S}'_{m_2} \otimes \boldsymbol{I}_K\right) & \boldsymbol{O} \\ \boldsymbol{O} & \boldsymbol{I}_{K(K+1)/2} \end{bmatrix} \begin{bmatrix} \boldsymbol{R} \\ \boldsymbol{S} \end{bmatrix}. \tag{A.12}$$

**Proof of the Auxiliary Result.** The proof uses mathematical induction so let us compare the impulse responses of the VAR($p$) and VAR($p+1$) models. For $p = 0$ we are comparing a white noise model with a VAR(1) with coefficient matrix $\boldsymbol{A}$. The impulse responses at horizon $h$ (the case $h = 1$ is trivial so we focus on $h > 1$) for these models are $\boldsymbol{O}_{K \times K}$ and $\boldsymbol{A}^h$, respectively. The base case $p = 0$ holds.

We start the inductive step by defining the companion matrix of the VAR($p+1$),

$$\boldsymbol{F}_{(p+1)} = \begin{bmatrix} \boldsymbol{A}_1 & \boldsymbol{A}_2 & \dots & \boldsymbol{A}_p & \boldsymbol{A}_{p+1} \\ \boldsymbol{I}_K & \boldsymbol{O} & \dots & \boldsymbol{O} & \boldsymbol{O} \\ \boldsymbol{O} & \boldsymbol{I}_K & \dots & \boldsymbol{O} & \boldsymbol{O} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \boldsymbol{O} & \boldsymbol{O} & \dots & \boldsymbol{I}_K & \boldsymbol{O} \end{bmatrix}. \tag{A.13}$$

This companion matrix is $(K(p+1) \times K(p+1))$. For an arbitrary matrix of this size, say $A$, we will introduce the notation $[A]^{ij}$ to denote its $(i,j)$th block of dimension $(K \times K)$. In this notation, the impulse response at horizon $h$ for the VAR$(p+1)$ is simply $[(F_{p+1})^h]^{11}$. Setting $A_{p+1} = O$ provides $K$ zero columns, and hence

$$[(F_{p+1})^h]^{11} = \sum_{j_1=1}^{p+1} \cdots \sum_{j_{h-1}=1}^{p+1} [F_{p+1}]^{1j_1}[F_{p+1}]^{j_1 j_2} \cdots [F_{p+1}]^{j_{h-1}1}$$

$$= \sum_{j_1=1}^{p} [F_{p+1}]^{1j_1} \left( \sum_{j_2=1}^{p+1} \cdots \sum_{j_{h-1}=1}^{p+1} [F_{p+1}]^{j_1 j_2} \cdots [F_{p+1}]^{j_{h-1}1} \right)$$

$$= \ldots = \sum_{j_1=1}^{p} \cdots \sum_{j_{h-1}=1}^{p} [F_{p+1}]^{1j_1}[F_{p+1}]^{j_1 j_2} \cdots [F_{p+1}]^{j_{h-1}1} = \left[ (F_p)^h \right]^{11}. \quad (A.14)$$

where $F_p$ is the companion matrix related to the VAR$(p)$. The inductive step is also complete.

*Proof of Theorem 2.* We rewrite the expression in Theorem 2 as

$$\sqrt{T} \left( \boldsymbol{\mu}(\hat{\boldsymbol{\theta}}_{T,m}, \mathbf{0}_{K^2(p-m)}, \hat{\boldsymbol{\sigma}}) - \boldsymbol{\mu}(\boldsymbol{\theta}_{T,p}, \boldsymbol{\sigma}) \right)$$

$$= \sqrt{T} \left( \boldsymbol{\mu}(\hat{\boldsymbol{\theta}}_{T,m}, \mathbf{0}_{K^2(p-m)}, \hat{\boldsymbol{\sigma}}) - \boldsymbol{\mu}(\boldsymbol{\theta}_{T,m}, \mathbf{0}_{K^2(p-m)}, \boldsymbol{\sigma}) \right)$$

$$- \sqrt{T} \left( \boldsymbol{\mu}(\boldsymbol{\theta}_{T,p}, \boldsymbol{\sigma}) - \boldsymbol{\mu}(\boldsymbol{\theta}_{T,m}, \mathbf{0}_{K^2(p-m)}, \boldsymbol{\sigma}) \right). \quad (A.15)$$

The first term in the RHS of Equation (A.15) contains a parameter transformation of the estimated parameters. The first-order delta method can be applied to this expression because Theorem 2 explicitly assumes the nonvanishing derivatives at the necessary points. The second term is nonrandom. It is the difference of two terms which only differ in locally misspecified coefficients which have been set to zero. A Taylor expansion can handle this second contribution. The result from the delta method is

$$\sqrt{T} \left( \boldsymbol{\mu}(\hat{\boldsymbol{\theta}}_{T,m}, \mathbf{0}_{K^2(p-m)}, \hat{\boldsymbol{\sigma}}) - \boldsymbol{\mu}(\boldsymbol{\theta}_{T,m}, \mathbf{0}_{K^2(p-m)}, \boldsymbol{\sigma}) \right)$$

$$\xrightarrow{d} \left( \frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta}_\infty, \boldsymbol{\sigma})}{\partial \boldsymbol{\theta}'} (\boldsymbol{S}_m \otimes \boldsymbol{I}_K) \right) \left( \boldsymbol{A}_m \boldsymbol{\delta} + \left( [\boldsymbol{S}_m' \boldsymbol{\Omega} \boldsymbol{S}_m]^{-1} \boldsymbol{S}_m' \otimes \boldsymbol{I}_K \right) \boldsymbol{R} \right) + \left( \frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta}_\infty, \boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}'} \right) \boldsymbol{S}, \quad (A.16)$$

where $\boldsymbol{\theta}_\infty$ denotes the parameter vector $\boldsymbol{\theta}_{T,p}$ but with all the misspecification parameters $\boldsymbol{C}_T$ set to zero. The result of the Taylor expansion is

$$\sqrt{T} \left( \boldsymbol{\mu}(\boldsymbol{\theta}_{T,p}, \boldsymbol{\sigma}) - \boldsymbol{\mu}(\boldsymbol{\theta}_{T,m}, \mathbf{0}_{K^2(p-m)}, \boldsymbol{\sigma}) \right)$$

$$= \sqrt{T} \boldsymbol{\mu}(\boldsymbol{\theta}_{T,m}, \mathbf{0}_{K^2(p-m)}, \boldsymbol{\sigma}) + \left( \frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta}_\infty, \boldsymbol{\sigma})}{\partial \boldsymbol{\theta}'} (\boldsymbol{S}_0 \otimes \boldsymbol{I}_K) \right) \left( (\boldsymbol{I}_{Kp_2} - \boldsymbol{\Pi}_m' \boldsymbol{\Pi}_m) \otimes \boldsymbol{I}_K \right) \boldsymbol{\delta} + O(T^{-1/2})$$

$$- \sqrt{T} \boldsymbol{\mu}(\boldsymbol{\theta}_{T,m}, \mathbf{0}_{K^2(p-m)}, \boldsymbol{\sigma})$$

$$= \frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta}_\infty, \boldsymbol{\sigma})}{\partial \boldsymbol{\theta}'} \left( \boldsymbol{S}_0 (\boldsymbol{I}_{Kp_2} - \boldsymbol{\Pi}_m' \boldsymbol{\Pi}_m) \otimes \boldsymbol{I}_K \right) \boldsymbol{\delta} + O(T^{-1/2}). \quad (A.17)$$

The notation can be made a little lighter using the definitions in Theorem 2: $D_\theta = \partial\mu(\theta_\infty, \sigma)/\partial\theta'$ and $D_\sigma = \partial\mu(\theta_\infty, \sigma)/\partial\sigma'$. Equations (A.15), (A.16) and (A.17) combine to

$$\sqrt{T}\Big(\mu(\hat{\theta}_{T,m}, \mathbf{0}_{p-m}, \hat{\sigma}) - \mu(\theta_{T,p}, \sigma)\Big)$$

$$\xrightarrow{d} D_\theta \Big[\big(S_m \otimes I_K\big) A_m - S_0 \big(I_{Kp_2} - \Pi'_m \Pi_m\big) \otimes I_K\Big] \delta$$

$$+ D_\theta \Big(S_m \big[S'_m \Omega S_m\big]^{-1} S'_m \otimes I_K\Big) R + D_\sigma S$$

$$= D_\theta \Big[\big(S_m \big[S'_m \Omega S_m\big]^{-1} S'_m \Omega - I_{Kp}\big) S_0 \big(I_{Kp_2} - \Pi'_m \Pi_m\big) \otimes I_K\Big] \delta$$

$$+ D_\theta \Big(S_m \big[S'_m \Omega S_m\big]^{-1} S'_m \otimes I_K\Big) R + D_\sigma S, \tag{A.18}$$

where the final line follows from the definition of $A_m$. With the given definitions of $C_m$ and $P_m$ we indeed recover the result stated in Theorem 2,

$$\sqrt{T}\Big(\mu(\hat{\theta}_{T,m}, \mathbf{0}_{K^2(p-m)}, \hat{\sigma}) - \mu(\theta_{T,p}, \sigma)\Big) \xrightarrow{d} D_\theta C_m \delta + D_\theta P_m R + D_\sigma S$$

$$\sim \ \mathrm{N}\Big(D_\theta C_m \delta, D_\theta P_m (\Omega \otimes \Sigma) P'_m D'_\theta + D_\sigma \Xi_{22} D'_\sigma\Big). \tag{A.19}$$

*Proof of Theorem 3.* By Equation (A.19),

$$\sqrt{T}\Big(\bar{\mu}(w) - \mu(\theta_{T,p}, \sigma)\Big) \ = \ \sum_{m=p_1}^{p} w_m \left[\sqrt{T}\Big(\mu(\hat{\theta}_{T,m}, \mathbf{0}_{K^2(p-m)}, \hat{\sigma}) - \mu(\theta_{T,p}, \sigma)\Big)\right]$$

$$\xrightarrow{d} D_\theta \sum_{m=p_1}^{p} w_m C_m \delta + D_\theta \sum_{m=p_1}^{p} w_m P_m R + D_\sigma S. \tag{A.20}$$

The calculation of the mean vector and the asymptotic covariance matrix is straightforward.

*Proof of Theorem 4.* A valid confidence interval for a scalar focus was derived in Theorem 6 of Liu (2015). We follow the same reasoning. Similar to Equation (A.20), we have

$$\sqrt{T}\Big(\bar{\mu}(\hat{w}) - \mu(\theta_T, \sigma)\Big) \xrightarrow{d} D_\theta \sum_{m=p_1}^{p} w_m(R_\delta) C_m R_\delta + D_\theta \big(\Omega^{-1} \otimes I_K\big) R + D_\sigma S. \tag{A.21}$$

Next, by the convergence of $\hat{\delta}$ to $R_\delta$,

$$\sqrt{T}\Big(\bar{\mu}(\hat{w}) - \mu(\theta_T, \sigma)\Big) - D_\theta \sum_{m=p_1}^{p} \hat{w}_m \hat{C}_m \hat{\delta} \xrightarrow{d} \mathrm{N}\big(\mathbf{0}, D_\theta(\Omega^{-1} \otimes \Sigma)D'_\theta + D_\sigma \Xi_{22} D'_\sigma\big). \tag{A.22}$$

The confidence region is constructed from the standardized quadratic form with population quantities replaced by their consistent estimates.

*Proof of Theorem 5.* Consider $m \neq p$. We define $\alpha_m = C'_m D'_\theta$, $\mathcal{A} = S'_0 \Omega^{-1} S_0 \otimes \Sigma$, and introduce a standard normally distributed random vector $Z_{K^2 p_2} \sim \mathrm{N}(\mathbf{0}, I_{K^2 p_2})$ and random variable $Z \sim \mathrm{N}(0, 1)$.

Now $R_\delta = \delta + \mathcal{A}^{1/2} Z_{K^2 p_2}$, and

$$
\begin{aligned}
R'_\delta C_m D'_\theta D_\theta C_m R_\delta &= \left(\alpha'_m(\delta + \mathcal{A}^{1/2} Z_{K^2 p_2})\right)^2 = \|\alpha'_m \mathcal{A}^{1/2}\|^2 \left(\frac{\alpha'_m \mathcal{A}^{1/2} Z_{K^2 p_2}}{\|\alpha'_m \mathcal{A}^{1/2}\|} + \frac{\alpha'_m \delta}{\|\alpha'_m \mathcal{A}^{1/2}\|}\right)^2 \\
&= (\alpha'_m \mathcal{A} \alpha_m) \left(Z + \alpha'_m \delta / \sqrt{\alpha'_m \mathcal{A} \alpha_m}\right)^2,
\end{aligned} \tag{A.23}
$$

where $\alpha'_m \mathcal{A} \alpha_m = a_m$. The squared expression has the stated noncentral chi-squared distribution (see Chapter 29 of Johnson et al. (1994) for details and moments). Finally, all quantities in $\widehat{FIC}_p = D_\theta(\hat{\Omega}^{-1} \otimes \hat{\Sigma}) D'_\theta + D_\sigma \Xi_{22} D'_\sigma$ are estimated consistently.

*Proof of Theorem 6.* If either $m = p$ and/or $l = p$, then there is no bias contribution and the matrix elements converge in probability. Now consider $m, l \neq p$, then
(a) For $m = l$, the proof is identical to the proof of Theorem 5, hence omitted.
(b) Start by noting that $x' A x = x'(\frac{A + A'}{2})x$, this gives

$$
R'_\delta \alpha_m \alpha'_l R_\delta = \left(\mathcal{A}^{-1/2} \delta + Z_{K^2 p_2}\right)' \left[\mathcal{A}^{1/2} \left(\frac{\alpha_m \alpha'_l + \alpha_l \alpha'_m}{2}\right) \mathcal{A}^{1/2}\right] \left(\mathcal{A}^{-1/2} \delta + Z_{K^2 p_2}\right) \tag{A.24}
$$

We subsequently use the transformation stated in Imhof (1961). The matrix in square brackets is symmetric and has a rank of at most two. The eigenvalue decomposition mentioned in Theorem 6 applies, and therefore

$$
R'_\delta \alpha_m \alpha'_l R_\delta = \sum_{i=1}^{2} \lambda_i \left(v'_i \mathcal{A}^{-1/2} \delta + Z_i\right)^2 \sim \sum_{i=1}^{2} \lambda_i \chi^2_{noncentral} \left(1, \left(v'_i \mathcal{A}^{-1/2} \delta\right)^2\right), \tag{A.25}
$$

where the independence of the $Z_i$ follows from orthonormality of the eigenvectors.

# Funding

# References

Abadir, K., Magnus, J. (2002). Notation in econometrics: A proposal for a standard. *The Econometrics Journal* 5:76–90.
Akaike, H. (1998). Information theory and an extension of the maximum likelihood principle. In: *Selected Papers of Hirotugu Akaike*. New York: Springer-Verlag, pp. 199–213.
Bates, J., Granger, C. (1969). The combination of forecasts. *Journal of the Operational Research Society* 20:451–468.
Benkwitz, A., Neumann, M., Lütkepohl, H. (2000). Problems related to confidence intervals for impulse responses of autoregressive processes. *Econometric Reviews* 19:69–103.
Berk, K. (1974). Consistent autoregressive spectral estimates. *The Annals of Statistics* 2:489–502.
Bruder, S. Wolf, M. (2017). Balanced bootstrap joint confidence bands for structural impulse response functions. *Working Paper*.
Burnham, K. P. Anderson, D. R. (2002). *Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach*. New York: Springer-Verlag.
Charkhi, A., Claeskens, G., Hansen, B. (2016). Minimum mean squared error model averaging in likelihood models. *Statistica Sinica* 26:809–840.
Claeskens, G., Croux, C., van Kerckhoven, J. (2007). Prediction-focused model selection for autoregressive models. *Australian & New Zealand Journal of Statistics* 49:359–379.
Claeskens, G., Hjort, N. L. (2003). The focused information criterion. *Journal of the American Statistical Association* 98:900–916.
Claeskens, G., Hjort, N. L. (2008). *Model Selection and Model Averaging*. Cambridge: Cambridge University Press.
Davidson, J. (1994). *Stochastic Limit Theory*. Oxford: Oxford University Press.

DiTraglia, F. (2016). Using invalid instruments on purpose: Focused moment selection and averaging for GMM. *Journal of Econometrics* 195:187–208.

Fujiwara, M. (1916). Ueber die obere Schranke des absoluten Betrages der Wurzeln einer algebraischen Gleichung. *Tohoku Mathematica Journal* 10:161–171.

Hamilton, J. D. (1994). *Time Series Analysis*. Princeton: Princeton University Press.

Hannan, E. (1970). *Multiple Time Series*. New York: John Wiley & Sons.

Hansen, B. E. (2005). Challenges for econometric model selection. *Econometric Theory* 21:60–68.

Hansen, B. E. (2007). Least squares model averaging. *Econometrica* 75:1175–1189.

Hansen, B. E. (2008). Least squares forecast averaging. *Journal of Econometrics* 146:342–350.

Hansen, B. E. (2016). Stein combination shrinkage for vector autoregressions. *Working Paper*.

Hansen, B. E., Racine, J. (2012). Jackknife model averaging. *Journal of Econometrics* 167:38–46.

Imhof, J.-P. (1961). Computing the distribution of quadratic forms in normal variables. *Biometrika* 48:419–426.

Ivanov, V., Kilian, L. (2005). A practitioner's guide to lag order selection for VAR impulse response analysis. *Studies in Nonlinear Dynamics & Econometrics* 9:1–36.

Johnson, N., Kotz, S., Balakrishnan, S. (1994). *Continuous Univariate Distributions, Vols. 1–2*. New York: John Wiley & Sons.

Kilian, L. (1998). Small-sample confidence intervals for impulse response functions. *Review of Economics and Statistics* 80:218–230.

Kilian, L. (2001). Impulse response analysis in vector autoregressions with unknown lag order. *Journal of Forecasting* 20:161–179.

Kitagawa, T., Muris, C. (2016). Model averaging in semiparametric estimation of treatment effects. *Journal of Econometrics* 193:271–289.

Lewis, R., Reinsel, G. (1985). Prediction of multivariate time series by autoregressive model fitting. *Journal of Multivariate Analysis* 16:393–411.

Liu, C.-A. (2015). Distribution theory of the least squares averaging estimator. *Journal of Econometrics* 186:142–159.

Liu, C.-A., Kuo, B.-S. (2016). Model averaging in predictive regressions. *Econometrics Journal* 19:203–231.

Liu, Q., Okui, R. (2013). Heteroscedasticity-robust Cp model averaging. *The Econometrics Journal* 16:463–472.

Lu, X. (2015). A covariate selection criterion for estimation of treatment effects. *Journal of Business & Economic Statistics* 33:506–522.

Lütkepohl, H. (1990). Asymptotic distributions of impulse response functions and forecast error variance decompositions of vector autoregressive models. *The Review of Economics and Statistics* 72:116–125.

Lütkepohl, H. (2005). *New Introduction to Multiple Time Series Analysis*. Berlin: Springer Science.

Lütkepohl, H., Staszewska-Bystrova, A., Winker, P. (2015). Comparison of methods for constructing joint confidence bands for impulse response functions. *International Journal of Forecasting* 31:782–798.

Rohan, N., Ramanathan, T. V. (2011). Order selection in arma models using the focused information criterion. *Australian & New Zealand Journal of Statistics* 53:217–231.

Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics* 6:461–464.

Zhang, X., Liu, C.-A. (2017). Inference after model averaging in linear regression models. *Working Paper*.

Zhang, X., Wan, A., Zou, G. (2013). Model averaging by jackknife criterion in models with dependent data. *Journal of Econometrics* 174:82–94.