

Sign Restrictions in Structural Vector Autoregressions: A Critical Review

RENÉE FRY AND ADRIAN PAGAN*

The paper provides a review of the estimation of structural vector autoregressions with sign restrictions. It is shown how sign restrictions solve the parametric identification problem present in structural systems but leaves the model identification problem unresolved. A market and a macro model are used to illustrate these points. Suggestions have been made on how to find a unique model. These are reviewed. An analysis is provided of whether one can recover the true impulse responses and what difficulties might arise when one wishes to use the impulse responses found with sign restrictions. (JEL C32, C51, E12)

1. Introduction

Structural vector autoregressions have become one of the major ways of extracting information about the macro economy. One might cite three major uses of them in macroeconometric research: for quantifying impulse responses to macroeconomic shocks; for measuring the degree of uncertainty about the impulse responses or other quantities formed from them; and for deciding on the contribution of different shocks to fluctuations and forecast errors through variance decompositions.

To determine this information, a vector autoregression (VAR) is first fitted to summarize the data and then a structural VAR (SVAR) is proposed whose structural equation errors are taken to be the economic shocks. The parameters of these structural equations are then estimated by utilizing the information in the VAR. The VAR is a reduced form that *summarizes* the data; the SVAR provides an *interpretation* of the data. As for any set of structural equations, recovery of the structural equation parameters (shocks) requires the use of identification restrictions that reduce the number of “free” parameters in the structural equations to the number that can be recovered from the information in the reduced form.

Five major methods for recovering the structural equation parameters (identifying the shocks) are present in the literature. Four of these explicitly utilize *parametric* restrictions. These involve the nature of the

*Fry: Australian National University. Pagan: University of Sydney. We are grateful to the editor and referees for some very helpful suggestions regarding the structure of this review. Fabio Canova and Matthias Paustian also gave us useful comments—the latter greatly clarifying the results of section 5.1. Pagan’s early work on this topic was supported by ESRC Grant 000 23-0244. Fry’s research was supported by ARC Grant #DP0664024.

structural equations. Parametric restrictions on these equations can vary according to whether particular variables appear in the latter (Cowles Commission), whether there is a recursive causal structure (Herman O. A. Wold 1951; Maurice H. Quenouille 1957; Christopher A. Sims 1980), and whether shocks have known short-run (Jordi Gali 1992) or long-run (Olivier Jean Blanchard and Danny Quah 1989) effects. In each case, the parametric restrictions free up enough instruments for the contemporaneous endogenous variables in the structural equations, thereby enabling the parameters of those equations to be estimated. Recently, a fifth method for estimating SVARs has arisen that employs *sign restrictions* upon the impulse responses as a way of identifying shocks (Jon Faust 1998; Harald Uhlig 2005; Fabio Canova and Gianni De Nicoló 2002). Applications of this method have been growing, as seen in the papers listed in table 1. The table is a subset of published studies and adopts a taxonomy that distinguishes between cases where only sign restrictions are used (often there are mixtures of sign and parametric restrictions), the type of shock (permanent or transitory), the number of shocks identified, and whether the sign restrictions come from a formal model or not. Consequently, it is worth examining this literature in more detail, and the aim of our paper is to exposit how the method works and to identify some of the difficulties that can arise in its application.

In practical work, it is often found that a combination of all the methods mentioned above need to be employed in order to be able to identify all the shocks of interest. We emphasize that which of the five methods mentioned above is used in practice does not depend on the data, but rather on the preferences of the investigator and/or of those who wish to utilize a SVAR to study some issue. These preferences may

well be incompatible as some users may feel that certain types of restrictions are more plausible than others. Prima facie it does seem likely that long-run and sign restrictions would be regarded as less restrictive than the other approaches, but without a specific context there can be no basis for recommending any particular approach. Each has difficulties and these need to be understood when making an informed judgment on their utility. Although it is likely in practice that a mixed set of restrictions will be employed, because the literature on sign restrictions is more recent than that of the other methods, it is convenient to simply assume that only sign restrictions are being employed.

Section 2 introduces the most common summative model (a VAR) and two structural models used in later analysis—a simple demand–supply model (called the market model) and a basic macroeconomic model determining output, interest rates, and inflation (called the macro model). Section 3 then examines how the five approaches described above would identify the shocks of the market model. Only a brief account of the parametric approaches is provided in order to allow for a more detailed description of the sign restriction methodology. By using the market model, it is shown that sign restrictions do implicitly impose parametric restrictions. In section 4, we outline some difficulties that can arise when implementing sign restrictions and the various solutions that have been proposed to them. One example is how one is to respond to the fact that a unique set of impulse responses is not available. This arises because, while the sign restrictions solve the *structural identification problem* by providing sufficient information to identify the structural parameters, they leave unresolved what Alan J. Preston (1978) called the *model identification problem*—the latter referring to the fact that there are many models with

TABLE 1
SUMMARY OF EMPIRICAL SVAR STUDIES EMPLOYING SIGN RESTRICTIONS

Fluctuations	Peersman (2005) STNI Rüffer/Sanchez/Shen (2007) STNI Sanchez (2007) STNF
Exchange rate	An (2006) STOI Farrant/Peersman (2006) STNF Lewis (2007) STNF Bjørnland/Halvorsen (2008) MTNI Scholl/Uhlig (2008) STNI
Fiscal policy	Mountford/Uhlig (2005, 2008) STNI Dungey/Fry (2009) MPTNI
Housing	Jarociński/Smets (2008) MTNI Vargas-Silva (2008) STOI
Monetary policy	Faust (1998) STOI Canova/De Nicoló (2002) STOF Mountford (2005) STNI Uhlig (2005) STOI Rafiq/Mallick (2008) STOI Scholl/Uhlig (2008) STNI
Technology	Francis/Owyang/Theodorou (2003) MPTOI Francis/Owyang/Roush (2005) MPTOF Dedola/Neri (2006) SPTOF Chari/Kehoe/McGrattan (2008) MPTNF Peersman/Straub (2009) STNF
Various	Hau/Rey (2004) STNF Eickmeier/Hofmann/Worms (2009) STNI Fujita (2011) STOI

Notes: Restriction Type: S=Sign only, M=Mixed

Shock Types: P=Permanent, T=Transitory

Number of Shocks: O=One only, N=Numerous

Restriction Source: F=Formal, I=Informal

identified parameters that provide the same fit to the data. Another is what one does if only the effects of a single shock, such as technology or money, is of interest? Section 5 addresses a range of more complex questions involving whether the methodology can recover the correct impulse responses and how one is to handle both permanent and transitory shocks. Finally, section 6 concludes.

2. *The VAR Representation and Two Simple Structural Models*

2.1 *VAR and SVAR Representations*

Most of the literature we deal with assumes that the data can be represented by a VAR (for simplicity we will make it of first order)

$$(1) \quad z_t = A_1 z_{t-1} + e_t,$$

where z_t is an $n \times 1$ vector of variables and e_t is a set of errors that have zero expectation, constant covariance matrix Ω and no serial correlation. From this, an interpretation of the data is provided through a SVAR

$$(2) \quad B_0 z_t = B_1 z_{t-1} + \varepsilon_t,$$

where ε_t are shocks that have zero mean, no serial correlation, constant variances, and no correlation between the individual shocks, i.e., $E(\varepsilon_{it}\varepsilon_{jt}) = 0$. Comparing (1) and (2) gives $B_0 e_t = \varepsilon_t$ i.e., the structural shocks ε_t we seek to measure are linear combinations of the VAR errors e_t . The latter can be estimated by the VAR residuals \hat{e}_t . To estimate the structural (economic) shocks, ε_t , then requires that one construct an appropriate set of weights (\hat{B}_0) on \hat{e}_t . Clearly the VAR is the reduced form of the structure set out in the SVAR.

The solution to the VAR(1) is the moving average (MA) form

$$(3) \quad z_t = D_0 e_t + D_1 e_{t-1} + D_2 e_{t-2} + \dots,$$

where D_j is the j th period impulse response of z_{t+j} to a unit change in e_t ($D_0 = I_n$). It follows that the MA form for the SVAR is

$$z_t = C_0 \varepsilon_t + C_1 \varepsilon_{t-1} \dots,$$

with the j th period impulse response of z_{t+j} to ε_t being $\hat{C}_j = D_j B_0^{-1} = D_j C_0$ as $C_0 = B_0^{-1}$. It is important to note that, since A_1 can be estimated by regressing z_t on z_{t-1} , and so does not require a structural model specification, D_j can be estimated from that information (in the first order case $D_j = A_1^j$). Hence, if one knows C_0 , one can find all the C_j without stipulating a structural model. For this reason, we will sometimes set $A_1 = 0$ in our illustrations of the various approaches, as that facilitates a focus upon how C_0 is

determined by each of them. Moreover, a failure to accurately estimate C_0 will mean that further C_j will be estimated inaccurately.

2.2 Two Simple Structural Models

2.2.1 A Market (Demand/Supply) Model

We take the case of a simple model comprising a demand and a supply function with associated shocks. This will be termed the market model. Specifically the SVAR system will be

$$(4) \quad q_t = -\beta p_t + \phi_{qq} q_{t-1} + \phi_{qp} p_{t-1} + \varepsilon_{Dt}$$

$$(5) \quad p_t = \gamma q_t + \phi_{pq} q_{t-1} + \phi_{pp} p_{t-1} + \varepsilon_{St},$$

where the shocks have expected values of zero and are assumed uncorrelated with standard deviations of σ_D and σ_S respectively. Hence, in terms of the SVAR discussion above, $z_t = \begin{bmatrix} q_t \\ p_t \end{bmatrix}$. The reduced form of the market model is a VAR(1) with the form

$$(6) \quad q_t = a_{qq} q_{t-1} + a_{qp} p_{t-1} + e_{1t}$$

$$(7) \quad p_t = a_{pq} q_{t-1} + a_{pp} p_{t-1} + e_{2t}.$$

Since the equations (4) and (5) are essentially identical for arbitrary parameter values, at this point there is nothing that distinguishes the demand (ε_{Dt}) and cost shocks (ε_{St}), and the task is to introduce extra information that does enable us to identify these. It would seem likely that most researchers would agree with the sign information in table 2 for the impact of positive shocks upon the contemporaneous variables (a positive movement in ε_{St} will mean a negative supply

TABLE 2
SIGN RESTRICTIONS FOR MARKET MODEL SHOCKS

Variable\shock	Model (demand/supply) shocks	
	Demand	Supply
p_t	+	+
q_t	+	–

side shock).¹ Since the patterns are distinct, this suggests that we are likely to be able to identify separate shocks. Indeed it is clearly going to be a requirement that shocks have distinct sign patterns in their effects on variables if we are to isolate them separately.

2.2.2 A Small Macro Model

A small macro model that is used a lot involves an output gap (y_t), inflation (π_t), and a policy interest rate (i_t). In terms of (1), the system variables are $z'_t = [y_t \ \pi_t \ i_t]$. A first-order SVAR model for these variables would then be

$$(8) \quad y_t = z'_{t-1} \gamma_y + \beta_{yi} i_t + \beta_{y\pi} \pi_t + \varepsilon_{yt}$$

$$(9) \quad \pi_t = z'_{t-1} \gamma_\pi + \beta_{\pi i} i_t + \beta_{\pi y} y_t + \varepsilon_{\pi t}$$

$$(10) \quad i_t = z'_{t-1} \gamma_i + \beta_{iy} y_t + \beta_{i\pi} \pi_t + \varepsilon_{it}.$$

The three shocks will be monetary policy (ε_{it}), a demand shock (ε_{yt}), and a cost-push (supply) shock ($\varepsilon_{\pi t}$). For simplicity, the shocks will be treated as having no serial

correlation, so that the reduced form is a VAR(1) of the form

$$(11) \quad y_t = z'_{t-1} \alpha_y + e_{1t}$$

$$(12) \quad \pi_t = z'_{t-1} \alpha_\pi + e_{2t}$$

$$(13) \quad i_t = z'_{t-1} \alpha_i + e_{3t}.$$

The signs of the contemporaneous effects to positive shocks will most likely be those in table 3. Again these are distinct and this enables the separation of the three structural shocks.

3. Identifying Shocks

3.1 The Parametric Approaches

In order to contrast the sign restriction approach to other methods of identifying shocks, let us think about how one might estimate the market model using the types of parametric restrictions distinguished in the introduction (we ignore the first possibility of constraining some coefficients of lagged values to zero). These restrictions are designed to identify the structural equations and hence the shocks.

- (a) If the system is assumed to be recursive, e.g., β is set to zero, then ordinary least squares (OLS) can be applied to the supply equation, since q_t is a function of ε_{Dt} and this is uncorrelated with

¹Although in later analyses we will always exhibit the pattern matrix in response to positive shocks, it needs to be recognized that a pattern for C_0 of $\begin{bmatrix} - & - \\ - & + \end{bmatrix}$ would also be consistent with demand and cost shocks, although with negative signs. So one needs to allow for this in any search. Of course $\begin{bmatrix} - & + \\ - & - \end{bmatrix}$ and $\begin{bmatrix} + & + \\ + & - \end{bmatrix}$ would also be acceptable. Clearly the need to check for the complete set of compatible sign restrictions will grow as the number of shocks increases.

TABLE 3
SIGN RESTRICTIONS FOR MACRO MODEL SHOCKS

Variable\shock	Demand	Cost-push	Interest rate
y_t	+	—	—
τ_t	+	+	—
i_t	+	+	+

ε_{st} . Three unknown parameters are left and there are three pieces of information to estimate them with—the estimated variances of p_t , q_t and the covariance of p_t and q_t .

- (b) A restriction that (say) a demand shock has no long-run effect upon the price would imply that $\phi_{pq} = -\gamma$, and so the supply curve would become a function of Δq_t and p_{t-1} . This implies that there is one less structural parameter to estimate in the supply curve and q_{t-1} is then freed up to be used as an instrument for Δq_t . Once the supply equation is estimated the demand equation can be found by using the residuals $\hat{\varepsilon}_{st}$ as an instrument for p_t .

- (c) An assumption that the short run effect of a demand shock upon prices is zero would imply that $\frac{\gamma}{(1+\beta\gamma)} = 0$ in the reduced form (VAR) equation for p_t

$$(14) \quad p_t = \psi_1 q_{t-1} + \psi_2 p_{t-1} + \frac{\gamma \varepsilon_{Dt}}{(1+\beta\gamma)} + \frac{\varepsilon_{st}}{(1+\beta\gamma)},$$

where ψ_j are functions of ϕ_{ij} and γ, β . Consequently, the VAR residual for p_t would not involve ε_{Dt} and so can be used as an instrument for p_t in the demand curve.

Thus, in all cases, the identification problem is solved by reducing the number of parameters to be estimated to three and by making available suitable instruments for estimation.

3.2 Sign Restrictions

Now a set of n estimated shocks \hat{e}_t will be available from the model we choose to be the summative one. For the market model, the VAR shocks e_{1t} and e_{2t} are from (6)–(7) while in the macro model the shocks e_{jt} are in (11)–(13). By combining them in an appropriate way, we can produce candidate structural shocks $\hat{\varepsilon}_t$ that are uncorrelated. Now there will be many such combinations. Some of them will produce impulse responses that have the correct signs, while others will not. Thus, in the market model case, there will only be some weights that produce shocks that respect the patterns in table 2. So our first task is to select an algorithm that gives a set of weights. Once one has these we can check if they are “successful” in the sense that the impulse response functions \hat{C}_j for the corresponding structural shocks agree with the postulated sign information. If they are not successful we will discard them and “draw” another set of weights.

Now the critical constraint needed in designing an algorithm to do this is that the generated weights must be such as to ensure that the constructed structural shocks $\hat{\varepsilon}_t$ are uncorrelated. Suppose we begin by first estimating a *recursive* VAR, e.g., in the market model we could act as if β was

zero. In that case, after estimation, we would have a set of shocks \hat{v}_t such that $\hat{e}_t = \hat{B}_0^{-1}\hat{v}_t$, where \hat{B}_0 is a lower triangular matrix, as this characterizes a recursive system.² By design, these shocks, \hat{v}_t , are uncorrelated. However, rather than work directly with such shocks, it is desirable to work with shocks that have unit variance, and this can be done by dividing each of the \hat{v}_{kt} by its standard deviation. Hence, let \hat{S} be the matrix that has the estimated standard deviations of the \hat{v}_t on the diagonal and zeros elsewhere. Then $\hat{e}_t = \hat{B}_0^{-1}\hat{S}\hat{S}^{-1}\hat{v}_t = \hat{T}\hat{\eta}_t$, where $\hat{\eta}_t = \hat{S}^{-1}\hat{v}_t$ are now regarded as structural shocks. These shocks possess unit variances and can be thought of as coming from a structural system $T^{-1}z_t = T^{-1}B_1z_{t-1} + \eta_t$. These $\hat{\eta}_t$ shocks will be termed our *base set*. Notice that they are just a rescaled version of the \hat{v}_t , so their nature has not changed.³

Now we form combinations of the $\hat{\eta}_t$ using a matrix Q , i.e., $\hat{\eta}_t^* = Q\hat{\eta}_t$. The $\hat{\eta}_t^*$ are candidates for “named” structural shocks, e.g. “supply” and “demand.” They need to be uncorrelated and so Q must be restricted. The appropriate restriction is that Q is a square matrix such that $Q'Q = QQ' = I_n$, since that means

$$\begin{aligned}\hat{e}_t &= \hat{T}Q'Q\hat{\eta}_t \\ &= \hat{T}^*\hat{\eta}_t^*,\end{aligned}$$

and $\text{cov}(\eta_t^* \eta_t^{*'}) = Q \text{cov}(\hat{\eta}_t \hat{\eta}_t') Q' = I_n$. Thus we have found a new set of shocks, $\hat{\eta}_t^*$, with the same covariance matrix as $\hat{\eta}_t$ (and which will

reproduce the $\text{var}(z_t)$), but which will have a different impact (\hat{T}^*) upon e_t and the variables z_t . It is this ability to create a large number of candidate shocks with varying impulse responses that is the basis of sign restriction methods. It is clearly very simple to construct all these shocks using programs that do matrix operations once we have a method for forming a Q with the property $Q'Q = QQ' = I_n$. There are many such Q s and we will refer to each as a “draw.”

How does one find a Q matrix? There are actually quite a few ways of doing this. The two most popular utilize Givens and Householder transformations (the latter is the basis of the QR decomposition used in many ill-conditioned regression problems), but this does not exhaust the possibilities. We provide an account of each of these and the relationship between them in the following subsections.

3.2.1 Givens Matrices

In the context of a three variable VAR (the macro model), a 3×3 Givens matrix Q_{12} has the form

$$Q_{12} = \begin{bmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

i.e., the matrix is the identity matrix in which the block consisting of the first and second columns and rows has been replaced by cosine and sine terms and θ lies between

²Although the residuals \hat{v}_t could be thought of as structural shocks, \hat{e}_t , we want to make the point that they are just shocks to begin the search process with, and there is no need to regard the recursive system as a plausible structure. It is the fact that shocks found from a recursive system are uncorrelated that makes them useful. Later we mention other ways of initiating the search.

³Numerically it is generally more efficient to estimate $\hat{\eta}_t$ by estimating the covariance matrix of the residuals $\hat{e}_t, \hat{\Omega}$, and then applying a Cholesky decomposition

$F^{-1}\hat{\Omega}F'^{-1} = I_n$ to form $\hat{\eta}_t = F^{-1}\hat{e}_t$ rather than first estimating a recursive system. The $\hat{\eta}_t$ constructed in this way will have unit variances and be uncorrelated. This is a useful way of proceeding since all that is needed to implement it is the estimated covariance matrix of the errors in the equations of the summative model. It also means that the summative model need not be a VAR. It could be a vector error correction model or a state space model and we use that fact in later sections.

0 and π .⁴ Q_{12} is called a Givens rotation. Then $Q'_{12}Q_{12} = I_3$ using the fact that $\cos^2\theta + \sin^2\theta = 1$. There are then three possible Givens rotations for a three variable system; the others being Q_{13} and Q_{23} . Each of the Q_{ij} depends on a separate parameter θ_k . In practice, most users of the approach have adopted the multiple of the basic set of Givens matrices as Q , e.g., in the three variable case we would use

$$Q_G(\theta) = Q_{12}(\theta_1) \times Q_{13}(\theta_2) \times Q_{23}(\theta_3).$$

It is clear that Q_G is orthogonal and so shocks formed as $\eta_t^* = Q_G \eta_t$ will be uncorrelated and their impact upon z_t will be $\hat{T}^* = \hat{T}Q'_G$.

Now, the matrix Q_G above depends upon three different θ_k . Canova and De Nicoló (2002) suggested that one make a grid of M values for each of the values of θ_k between 0 and π , and then compute all the possible Q_G . Of course all of these models distinguished by different numerical values for θ_k are observationally equivalent in that they produce an exact fit to the variance of the data on z_t .⁵ Only those Q_G producing shocks that agree with the maintained sign restrictions would be retained.

As an example, we look at the macro model described in Seonghoon Cho and Antonio Moreno (2006) estimated with some data on the U.S. output gap, inflation, and the Federal Funds rate. As described above, begin with a recursive model imposing $\beta_{yi} = 0, \beta_{y\pi} = 0, \beta_{\pi i} = 0$. OLS on each of (8)–(10) then gives structural equation residuals that are uncorrelated. More potential structural shocks can subsequently be found by combining these residuals (after rescaling to make them have unit variances) with Q matrices. Two of these Q matrices (from

the Givens approach) are given below. They have the property that $Q'Q = I_3$.⁶

$$(15) \quad Q^{(1)} = \begin{bmatrix} -.4551 & .3848 & .8030 \\ -.5853 & -.8089 & .0559 \\ .6710 & -.4446 & .5933 \end{bmatrix},$$

$$Q^{(2)} = \begin{bmatrix} .0444 & -.8431 & .5359 \\ .8612 & -.2395 & -.4482 \\ .5062 & .4815 & .7155 \end{bmatrix}.$$

When $Q^{(2)}$ is used, the generated structural shocks have a sign pattern for C_0 of

$$\begin{bmatrix} + & - & + \\ + & + & + \\ + & + & + \end{bmatrix}, \text{ which disagrees with the}$$

restrictions in table 3. In contrast, $Q^{(1)}$ does produce a set of impulse responses that is consistent with the table. Hence, employing the sign restriction methodology only the impulses found using $Q^{(1)}$ would be retained.

3.2.2 Householder Transformations

The alternative method of forming an orthogonal matrix Q is to generate some 3×3 random variables W from an $N(0, I_3)$ density (for a three variable VAR) and then decompose $W = Q_R R$, where Q_R is an orthogonal matrix and R is a triangular matrix. Householder transformations of a matrix are used to decompose W . The algorithm producing Q_R is often called a QR decomposition. Clearly $Q_R = I$ corresponds to the matrix used in recursive orderings. Since many draws of W can be made, one can find many Q_R . Juan Francisco Rubio-Ramírez, Daniel Waggoner, and Tao Zha (2005) seem to have been the first to propose this, and they have argued that, as the size of the VAR grows, this is a computationally efficient strategy relative to the Givens approach. In Renée Fry and Adrian Pagan

⁴In general, Q_{ij} is formed by taking a $n \times n$ identity matrix and setting $Q_{ij}^{\theta} = \cos \theta$, $Q_{ji}^{\theta} = -\sin \theta$, $Q_{ij}^{\theta} = \sin \theta$, $Q_{ji}^{\theta} = \cos \theta$, where the superscripts refer to the row and column of Q_{ij} .

⁵It is important in the analysis that the z_t have been mean corrected before the VAR is fitted.

⁶The fact that we retain only four decimal places above means that $Q'Q$ is not exactly I_3 .

(2007), we show that the methods are equivalent, so the main factor in choice would be computational speed. As the system grows in size, we would expect the Householder method to be superior.

3.3 Sign Restrictions in the Market Model

So how do sign restrictions resolve the structural identification problem in the market model? As noted previously, the key problem is how to identify the initial impulse responses C_0 . For illustrative purposes, it is convenient to suppress the dynamics by setting $B_1 = 0$ and to study the solutions for C_0 alone. Information on the signs expected for the elements of C_0 is given in table 2. In the next step, a recursive system is set up and estimated. Of course the market model is generally not recursive, but this is simply a mathematical device to generate a set of shocks \hat{v}_t that are uncorrelated. We therefore assume the following recursive system

$$(16) \quad p_t = v_{1t}$$

$$(17) \quad q_t - \tau p_t = v_{2t}.$$

Three parameters are estimated in this system— $\tau = \text{cov}(q_t, p_t) / \text{var}(p_t)$ and the two variances of v_{jt} , σ_j^2 . Effectively, this means that the variances of p_t and q_t and their covariance are used for estimation. A base set of impulses will then be found from $\eta_{jt} = \sigma^j v_{jt}$, where σ^j is the inverse of the standard deviation of v_{jt} . By definition, $\eta_{1t} = \sigma^1 v_{1t} = \sigma^1 p_t$ and $\eta_{2t} = \sigma^2 v_{2t} = \sigma^2(q_t - \tau p_t)$.

These base impulses are then used to construct new shocks η_{jt}^* by using a Givens rotation as the weighting matrix. Since there is only one Givens matrix in the two variable case, $Q = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$, the transformed system becomes

$$\begin{bmatrix} \eta_{1t}^* \\ \eta_{2t}^* \end{bmatrix} = \begin{bmatrix} \sigma^1 p_t \cos \theta - \sigma^2 (q_t - \tau p_t) \sin \theta \\ \sigma^1 p_t \sin \theta + \sigma^2 (q_t - \tau p_t) \cos \theta \end{bmatrix},$$

Letting $\phi_1 = \cos \theta$ and $\phi_2 = \sin \theta$ the two equations can be written as

$$(18) \quad (\sigma^1 \phi_1 + \sigma^2 \phi_2 \tau) p_t - \sigma^2 \phi_2 q_t = \eta_{1t}^*$$

$$(19) \quad (\sigma^1 \phi_2 - \sigma^2 \tau \phi_1) p_t + \sigma^2 \phi_1 q_t = \eta_{2t}^*,$$

with impulse responses of (p_t, q_t) to η_{jt}^* being

$$\begin{bmatrix} \sigma^1 \phi_1 + \sigma^2 \phi_2 \tau & -\sigma^2 \phi_2 \\ \sigma^1 \phi_2 - \sigma^2 \tau \phi_1 & \sigma^2 \phi_1 \end{bmatrix}^{-1} = G^{-1} \\ = \frac{1}{\det(G)} \begin{bmatrix} \sigma^2 \phi_1 & \sigma^2 \phi_2 \\ -\sigma^1 \phi_2 + \sigma^2 \tau \phi_1 & \sigma^1 \phi_1 + \sigma^2 \phi_2 \tau \end{bmatrix}.$$

Because σ^j are fixed by the data, the sign of the impact of the shocks upon p_t and q_t will be dependent on $\text{sgn}(\phi_1)$ and $\text{sgn}(\phi_2)$, and these can be positive or negative depending upon the values taken by θ . Consequently there may be many impulse responses which satisfy the sign restrictions, each of which is indexed by a value of θ . Note that, even though we started with a recursive system, we will generally not have one as θ varies.

It is useful now to observe that, given θ , the number of unknown parameters in (18)–(19) has been reduced to three (τ, σ^1, σ^2), just as happened with the parametric methods. This reduction means that a unique set of parameters can be recovered for a given system of equations (and a specified θ), and so the structural parameter identification problem has been solved.

4. Some Basic Issues with Sign Restrictions

4.1 Model and Structural Identification

Now in the discussion of the previous section we only retain those shocks whose impulses agreed with the postulated signs. But it is clear that there may be many impulse

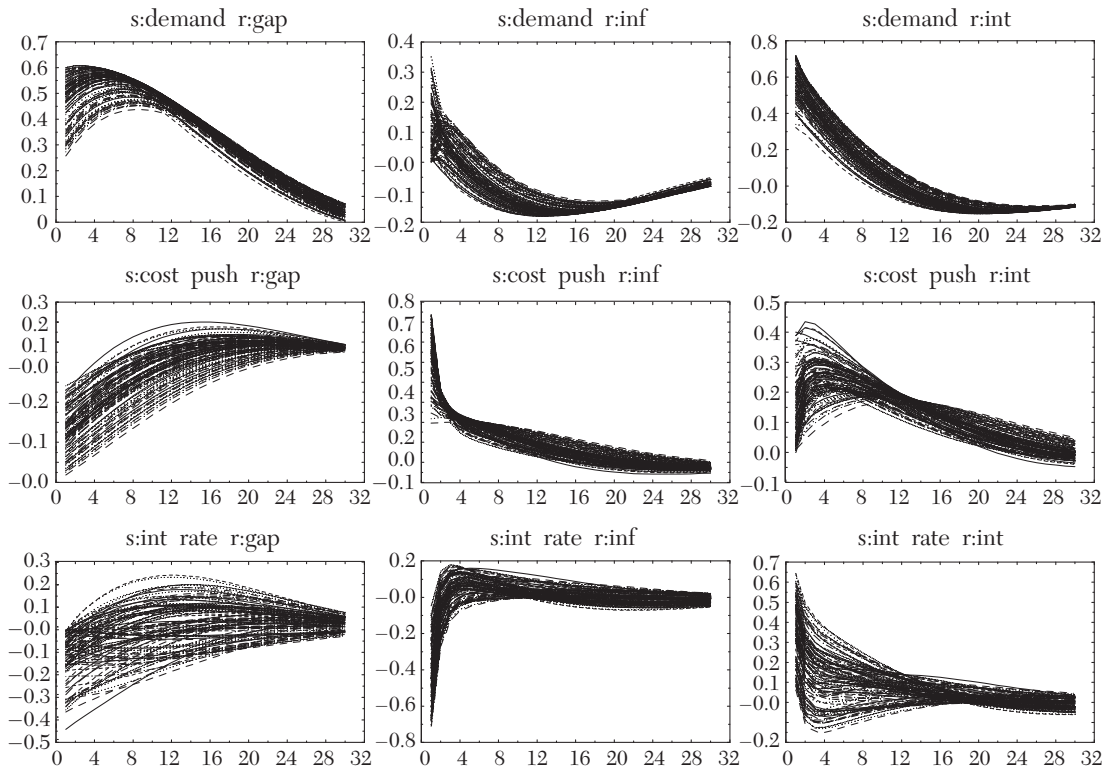


Figure 1. 1,000 Impulse Responses Found with the Sign Restrictions of Table 2

responses that satisfy these sign restrictions. Thus, when using Givens matrices, it is unlikely that there will be a single value of θ that will produce the requisite sign restrictions. Figure 1 shows the large range of impulse responses one gets by applying the contemporaneous sign restrictions of table 3 to the macro model data that we used earlier when illustrating the effects of choosing two values for Q . It is noticeable that, even though the initial response of output to the interest rate has been forced to be negative, there are some cases where that response becomes positive very quickly. Each value of θ produces a new *model* constituting a new set of structural equations and shocks. Consequently, although we have converted

any given system of equations (consistent with a given θ) to one that has a structure that is identified, we have not identified a unique *model*. The difference between structural and model identification was emphasized by Preston (1978).

What should one do about this *multiple models problem*? One response is to try to summarize the information in the graphs in some way, e.g., reporting a central tendency and the magnitude of the spread of responses. Thus, if the impulse responses $C_j^{(k)}$ that satisfy the sign restrictions are computed, where k indexes the different values of θ , various percentiles, such as the 5 percent, 50 percent, and 95 percent might be reported. This is done in figure 2.

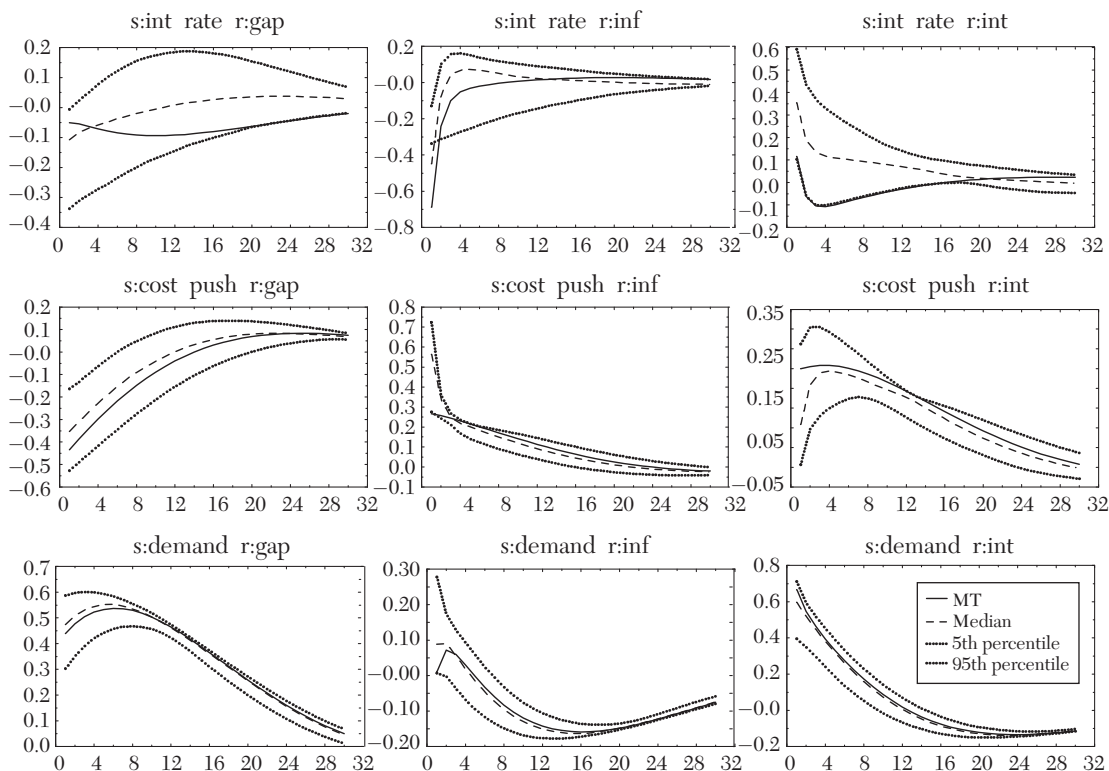


Figure 2. Impulse Responses for Macro Model: MT, Median, and 5, 95 Percentiles

It may seem as if this is emulating the approach when one presents percentiles of a distribution from either a Bayesian or bootstrap experiment. But it is important to recognize that the distribution here is *across models*. It has nothing to do with sampling uncertainty. Referring to this range as if it is a confidence interval (something that is very common in the applied literature) is quite false. All you get is a glimpse of the possible range of responses as the model varies. Of course, this might be valuable. Often we do present information about how our answers change as models are varied. Edward E. Leamer (1978) did this in a regression context with his extreme bounds analysis. But

it should not be imbued with probabilistic language. Even if the VAR parameters A_1 and Ω were known with certainty, there will be a question of how one proceeds whenever there are many θ . There is of course a greater range when one accounts for the uncertainty in A_1 and Ω , as is often done in this literature. Examples are Uhlig (2005) and Gert Peersman (2005), where Bayesian methods are applied to estimate the summative model, but it does not help to understand the model identification issue by confusing these two sources of variation.

Do any difficulties arise in interpreting (say) these summary measures? Let us illustrate the issues by considering the median

of the impulse responses, by far the most popular choice for capturing the central tendency. Suppose there is a single variable and two shocks, and that we have impulse responses $C_{11}^{(k)}$ and $C_{12}^{(k)}$, where k indexes the models (values for θ). Ordering these into ascending order enables one to find the medians $C_{11}^{(k_1)} = \text{med}\{C_{11}^{(k)}\}$ and $C_{12}^{(k_2)} = \text{med}\{C_{12}^{(k)}\}$. But k_1 may not equal k_2 , and so the model that produced the impulse response that is the median of $\{C_{11}^{(k)}\}$ may not be the same as that for $\{C_{12}^{(k)}\}$. Presenting the medians may be likened to presenting the responses to technology shocks from a real business cycle model, and the monetary shocks from a monetary model, and it is hard to believe that this is a reasonable approach. Clearly, this comment applies to other percentiles so that the extreme values that are being reported may come from very different models.

Another piece of information presented in many papers is a variance decomposition of $\text{var}(z_t)$ into the contributions from various shocks. Often this is done using the medians of the impulse responses. Is this correct? Now a variance decomposition requires that the shocks be defined in the same way, necessitating a common value of θ be used, since only in this case will the shocks be uncorrelated by design. If shocks are not uncorrelated then a variance decomposition does not make much sense. This issue is not resolved by another common practice that computes the fraction of the variance explained by the j th shock ($j = 1, \dots, n$) in the k th model ($k = 1, \dots, M$), $\psi_j^{(k)}$, $j = 1, \dots, n$, and then reports the n medians of $\{\psi_j^{(k)}\}_{k=1, \dots, M}$. Since in general these medians will not come from the same model, there is nothing that ensures that the $\text{med}\{\psi_j^{(k)}\}$ sum to one across all shocks, i.e., the variance is exhaustively accounted for.

Now it needs to be said that the issue of model identification is always present and is not specific to sign restrictions. Thus, if one

used a recursive system to get structural identification, there are many other such systems (orderings) that will yield the same VAR and give the same fit to the data. Each structure coming from a given ordering is parametrically identified but, as all of the orderings exactly replicate the data, there is no unique model. Only if one is prepared to consider that there is a single recursive model that is tenable as the data generating process will this occur, and that is rare. Indeed, one often sees comments to the effect that reordering the equations did not modify the conclusions much.

Why then should we pay any more attention to this model identification issue for sign restrictions than for other ways of identifying VARs? Some insight into this comes from examining the two possible recursive versions of the market model—that given in (16) and (17) and the other being where p_t and q_t are interchanged in these equations. Although observationally equivalent, the two models can be treated as different views about how the market operates. In one case, quantity is treated as predetermined, and so prices reconcile supply and demand, while the other has price being set and quantity doing the adjustment. A choice between these might be made using institutional knowledge that is difficult to put into a VAR framework. But, in the sign restriction approach to the market model, there is no equivalent interpretation, as the restriction employed for identification essentially ties the supply and demand elasticities together. Nevertheless, any solution to the multiple models problem has to be the same as for recursive models, namely the introduction of extra information that enables one to discriminate between them.

What sort of extra knowledge might be used? There is no one way of doing this in the literature. One possibility is to continue to add on sign restrictions that relate to longer lags in the impulses. Thus one can see

from the curves in figure 1 that, imposing a negative effect of an interest rate shock upon output and inflation for ten periods rather than one period, would eliminate many of the 1,000 models in that figure. To formally understand why this might work and the limitations to it, we examine the relation between impulses noted earlier viz. $C_j = D_j C_0$. Because D_j can be estimated from the VAR and does not need structural information, restrictions on C_j translate into *indirect* restrictions on C_0 , and so $C_j > 0$ requires C_0 to be such that $D_j C_0 > 0$. This may well reduce the number of possible C_0 's (models) that jointly satisfy sign restrictions on the higher order impulses as well as on C_0 . If the restrictions on C_j were quantitative, then the indirect restrictions implied on C_0 would certainly narrow the possible values of C_0 . However, the same effect does not necessarily hold for qualitative restrictions. For example, if all elements of C_0 are positive, and so too are the estimated D_1 from a VAR(1) fitted to the data, then a restriction that the elements of C_1 are positive adds nothing to what has already been assumed about the signs of C_0 . There appears to be a belief in the literature that adding on sign restrictions for longer impulse responses, $C_{j,j} > 0$, provides stronger identifying information, and this seems to stem from the Monte Carlo study in Matthias Paustian (2007). However, as is clear from the connections that exist between the C_j and C_0 noted above, nothing guarantees this.

Quantitative information about the likely magnitude of the impulse responses is sometimes invoked in order to reduce the set of models. Thus Kilian and Murphy (2009) argue that some estimates generated of the short-run supply elasticity of oil and the initial impact of oil prices upon activity are implausible, and so models producing them should be discarded. A second group of methods looks at setting up a criterion

based upon the magnitude of impulses and minimizing it with respect to θ . Faust (1998) and Uhlig (2005) do this. Uhlig's criterion is to give a high weight to "large" standardized impulses over "small" ones. Thus he says "Given a choice among many candidate monetary impulse vectors . . . it might therefore be desirable to pick the one which generates a more decisive response of the variables" (Uhlig 2005, p. 414). The exact form of the penalty function varies with the application. Thus, in general, all one can say is that a value for θ is found by minimizing a criterion that is a function of the magnitude of the impulse responses $C_j^{(k)}$. Provided it is clear that this is being done, it is simply a matter of deciding if the supplementary quantitative criterion is acceptable, although one needs to recognize that non-sign information is being invoked to get a unique model.

A different approach to selecting a single value of θ by minimizing a criterion is that in Fry and Pagan (2005). This begins with the observation that researchers seem to be attracted to the idea of presenting the median as a good summary of the central tendency of impulse responses across models. Our second observation was that the median responses may come from different models, potentially making them impossible to utilize in exercises such as variance decompositions. So it seems logical to find a single model whose impulse responses are as close to the median values as possible. We will term this the median target (MT) method. The MT solution is to choose that value of $\theta^{(k)}$ that produces impulses that are as close to the median responses as possible. To devise a criterion to do this, it is necessary to recognize that the impulses need to be made unit-free by standardizing them. This is done by subtracting off their median and dividing by their standard deviation, where these are measured over whatever set of models has been retained

as satisfying the sign restrictions. These standardized impulses are then placed in a vector $\phi^{(l)}$ (in a two variable case, ϕ is 4×1 as there are four impulses) for each value $\theta^{(l)}$. Subsequently we choose the l that minimizes $MT = \phi^{(l)'}\phi^{(l)}$, and then use that $\theta^{(l)}$ to calculate impulses. Whether this strategy produces a unique l is an empirical question, although in applications we have made it turns out to do so. In the event that the median shocks are uncorrelated, then we would find that the median responses would be selected by this criterion. So a difference between the median responses and the MT-selected responses essentially indicates that the shocks associated with the median impulses are correlated. Consequently, the MT methodology can be regarded as a diagnostic device.

Figure 2 also shows the median impulses and those coming from the MT approach. Clearly major differences in the effects of an interest rate shock upon output emerge when it is insisted that the shocks must come from a single model. A comparison of the median with the adjusted measure for other shocks does not reveal as great a difference, except perhaps in the initial impact of monetary policy on inflation. In Fry and Pagan (2005), we found that applying the MT method to the data in Blanchard and Quah (1989) produced very little difference when assessing the impact of demand and supply shocks. A number of other papers also report that the results are not too dissimilar, e.g., Rasmus Ruffer, Marcelo Sanchez, and Jian-Guang Shen (2007) and Canova and Paustian (2010). It does seem to us however that it is more satisfactory to ensure that the impulses come from the same model rather than getting them from different models, even if in some specific instances the adjustment does not produce major changes. At the very least, one needs to check that a failure to insist that shocks come from a single model has not created any distortions. The adjustment is

simple to compute. One might also observe that, although the discussion above has been about the median, it also applies to any of the “percentile” measures.

4.2 Identifying a Single Shock

In the description above, it was assumed that n shocks were to be found. But sometimes only a single shock is of interest and therefore only one shock is isolated. Examples are M. S. Rafiq and S. K. Mallick (2008), Lian An (2006), and Uhlig (2005), although there are many others where the number of shocks identified is greater than one but less than n . Dealing with the single shock case, we might still utilize the $n \times n$ Q -matrices above to produce n uncorrelated structural shocks, but only focus upon one of them. Two issues arise here. Firstly, in some papers one has the impression that it is only necessary that the weights used for constructing the structural shocks be a $n \times 1$ vector q that has unit length. Uhlig’s papers often state it in this way, e.g., Almuth Scholl and Uhlig (2008, p. 5). If q is not selected from a Q that is orthogonal, then the resulting shock need not be uncorrelated with the remaining (unidentified) $n - 1$ shocks. To the extent that one does not need this property for analysis, then there is no problem, but if one is trying to perform a variance decomposition it is mandatory. Our reading of a number of papers in the literature is that q was not selected in a way to ensure orthogonality.

A second problem arises from the following scenario. Suppose that there are two variables and we believe that one shock has a positive initial effect on the first variable. However we are unwilling to describe either its effects on the second variable or to set any signs for the initial effects of the second shock. This scenario would generate signs for C_0 of $\begin{bmatrix} + & ? \\ ? & ? \end{bmatrix}$, where $?$ means that no sign information is provided. It is clear

that this is not enough information to discriminate between the shocks. Indeed, even the pattern $\begin{bmatrix} + & ? \\ + & ? \end{bmatrix}$ would not suffice, since it is possible that the impulse responses found from a draw of Q might be $\begin{bmatrix} + & + \\ + & + \end{bmatrix}$, and then we are faced with the fact that both shocks have the same sign pattern. In any finite number of draws, one may not encounter this, but that is just fortuitous. Hence a problem arises if there is a failure to specify enough information to discriminate between shocks. We will refer to this as the *multiple shocks problem*, as distinct from the *multiple models problem* that was mentioned earlier.

As an illustration of the multiple shocks problem, suppose we look at the macro model when $Q^{(2)}$ in (weights) is used as the weighting matrix to form shocks. Suppose it is desired to identify only a single shock—demand—using the sign restrictions from the macro model. Then, when $Q^{(2)}$ is used to construct three shocks, two of these would produce the right signs, and so both are potential demand shocks. However, we cannot accept both as demand shocks given that they are in the same model. It is only if one describes the signs patterns for all of the shocks that it is possible to rule out the use of $Q^{(2)}$. Consequently, if only a single shock is to be isolated (more generally any number less than n) some information will need to be provided on what strategy was used to deal with this issue. At the moment little mention is made in many published articles using sign restrictions. The problem does seem to come up quite a bit, e.g., Rüffer, Sanchez, and Shen (2007) mention that it occurs in their study, although they offer no comment on what they did about it, and Andrew Mountford (2005) also alludes to it.

4.3 The Origin of Sign Restrictions

Generally these have been rather informal, although increasingly they have been drawn from dynamic stochastic general

equilibrium (DSGE) models. Thus the New Keynesian (NK) policy model with the form

$$\begin{aligned} y_t &= \alpha_{1y} y_{t-1} + \beta_{1y} E_t(y_{t+1}) \\ &\quad + \gamma_{1i}(i_t - E_t(\pi_{t+1})) + \varepsilon_{yt} \\ \pi_t &= \alpha_{2\pi} \pi_{t-1} + \beta_{2\pi} E_t(\pi_{t+1}) \\ &\quad + \gamma_{2y} y_t + \varepsilon_{\pi t} \\ i_t &= \alpha_{3\pi} i_{t-1} + \gamma_{3y} y_t + \beta_{3\pi} E_t \pi_{t+1} + \varepsilon_{it}, \end{aligned}$$

is often invoked as a small macro model. Assuming that there is no serial correlation in the shocks, the solution is a VAR(1) in $z'_t = [y_t \ \pi_t \ i_t]$, with the VAR(1) coefficients A_1, Ω being a function of the NK model parameters θ . Maximum likelihood estimates of θ can be found with the same data as used to fit the VAR(1) for the sign restriction work. Using the maximum likelihood estimates of θ , the resulting impulse responses are in figure 3, along with those coming from the MT method of producing a unique set of impulses under sign restrictions. There are some very large quantitative differences. Indeed, the NK impulses often lie well outside the range coming from the 1,000 models produced with sign restrictions. It might be wondered why this is the case as the NK model implies a VAR for the variables. But it is a restricted VAR. Hence, if the sign restriction impulses are $C_j^{SR} = C_0^{SR} \hat{D}_j^{VAR}$, then those from the NK model will be governed by $C_j^{NK} = C_0^{NK} \hat{D}_j^{NK}$. Consequently, there can be two reasons for the difference between impulses from the two approaches—a discrepancy between the initial effects C_0^{SR} and C_0^{NK} and a difference between the estimated VAR coefficients D_j .⁷

⁷In this case, the coefficients of π_{t-1} in the VAR equation for π_t are 0.4 (NK) and 0.9 (unrestricted). The covariance matrix of the VAR errors also shows some differences—the correlation between the output gap and inflation VAR equation errors being 0.4 in the data and 0.1 implied by the NK model.

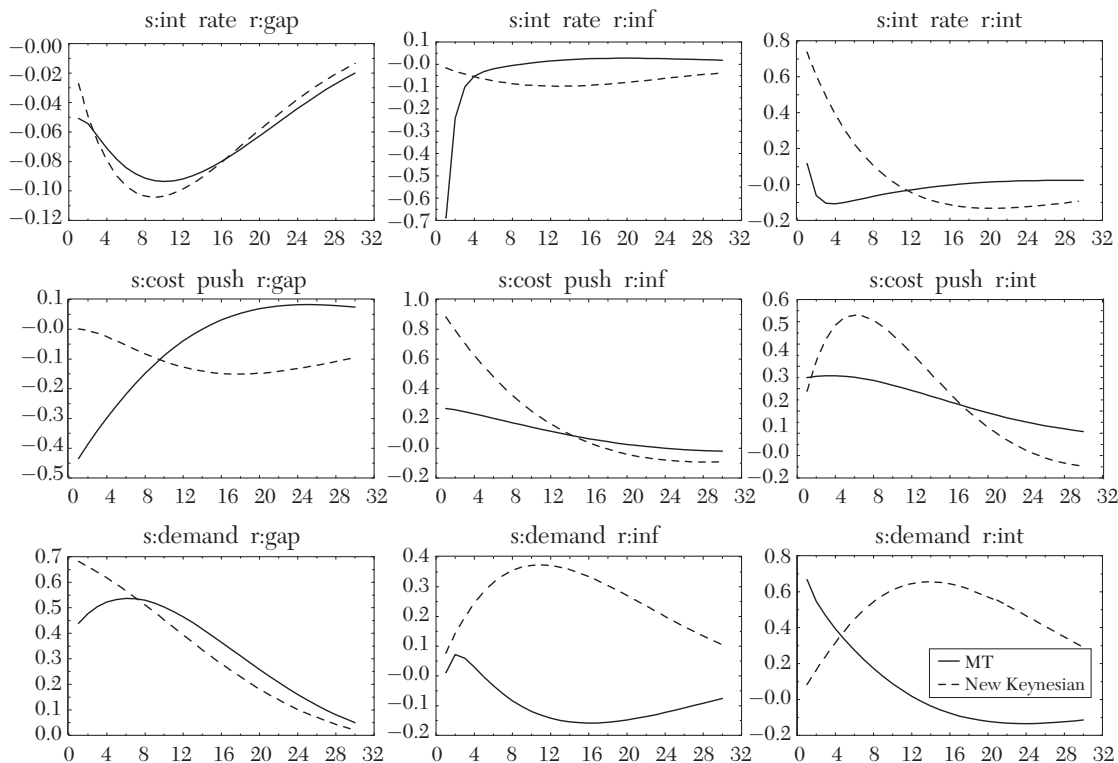


Figure 3. Comparison of the Impulse Responses of the Small Macro Model (MT method—solid line), and a New Keynesian Model (dashed line)

All the sign restriction models keep D_j fixed at the OLS estimates of the VAR(1). In contrast, the NK model says there are restrictions on the VAR parameters, and these are imposed in the maximum likelihood estimation. Hence, if the NK model restrictions are incorrect, there will be a bias in \hat{D}_j^{NK} , which will show up as different impulse responses to the unrestricted ones. This points to a rationale for just using the NK model as a source of sign restrictions rather than exploiting its stronger implications for VAR coefficient relations. Of course the sign restrictions may depend upon the structural model coefficient values fed into the NK model, and so it has been proposed that the model be simulated for a wide range of these, retaining only

those signs that are robust to the parameter variations. The sign restrictions in table 3 are likely to be broadly consistent with those found by simulating models, such as the NK above, for a range of parameter values.

As we will see in the next section, this strategy of using theoretic models to produce sign restrictions has been increasingly used in the literature. Canova and Paustian (2010) examine it in some detail, simulating data from a DSGE model and then seeing if the correct impulse responses would be recovered. They find that it recovers the shocks reasonably well, provided that enough of these are used and all shocks are identified. In contrast, Jarkko Jääskelä and Jennings (2010) found that they could not

recover the correct impulse responses from an NK model of a small open economy, despite using many sign restrictions.

The model-based approach to producing sign restrictions seems a useful way to proceed, as it does not commit the user to the DSGE model, but has the advantage that it restricts the informal approach in a fashion that probably commands reasonable assent. A lot depends on why one is performing the VAR analysis. If one is trying to “discover” what the data says about relations then imposing sign restrictions from (say) the NK model above would not appeal as much, since one would never see (say) that interest rates had a positive impact on inflation in the data. “Puzzles” like this are sometimes the source of productive theorizing and so one should be careful about predetermining outcomes. Of course one check on this is available from the draws that yielded impulse responses that *didn’t agree* with the sign restrictions. If there are a large number of these, then one might well conclude that the evidence is compatible with too many models to make one comfortable with the idea of restricting attention to those satisfying the sign restrictions. Sometimes the number rejected is very high, e.g., Kilian and Murphy (2009) report only 30,860 “successful” models from 1.5 million draws. Roland Straub and Peersman (2006) used the rejected information in this way to assess whether the NK model was a good description of the data.

5. *Some Advanced Issues with Sign Restrictions*

A number of questions and issues arise with sign restrictions that deserve comment. First, do sign restrictions recover the correct impulse responses? This question is considered in the next subsection by using the market model, and it is concluded that, while they can be potentially recovered up to an unknown scaling factor, the standard

strategies for dealing with the multiple models problem may mean that the true impulse responses are not isolated. But even if there is no certainty that the correct impulses can be found, it is desirable to maximize the chances of doing so and we therefore examine some recommendations that have been made about how to do this. Second, if we wish to align the summative model with theoretical models, it is often necessary to recognize that, whilst the latter generally imply a VAR in *all* the model variables, only a subset of variables are actually used in modeling, and these may not follow a VAR, i.e., using a VAR as the summative model would be incorrect. We show how this complication can be dealt with fairly easily. Finally, we ask what one does if there are both permanent and transitory shocks in the system? Again a VAR is not the correct summative model and it needs to be replaced by a VECM. Hence we spend some time indicating how to adapt the methods proposed earlier in connection with VARs to the vector error correction model (VECM) case.

5.1 *Can We Recover Correct Impulse Responses from Sign Restrictions?*

In the literature, one sometimes gains the impression that the answer to the question posed above is in the affirmative. But it is a tricky question to give an answer to, as it depends on what type of experiment you wish to perform with the impulse responses. To see why, note that, in the market model with no dynamics, the VAR equations for p_t and q_t would be

$$(20) \quad q_t = \left(1 - \frac{\beta\gamma}{(1 + \beta\gamma)}\right)\varepsilon_{Dt} - \frac{\beta\varepsilon_{St}}{(1 + \beta\gamma)}$$

$$(21) \quad p_t = \frac{\gamma\varepsilon_{Dt}}{(1 + \beta\gamma)} + \frac{\varepsilon_{St}}{(1 + \beta\gamma)}.$$

Now, because ε_{St} and ε_{Dt} are uncorrelated, dividing by their standard deviations would

produce some base shocks η_{jt} that have unit variance, namely $\eta_{1t} = \sigma_D^{-1}\varepsilon_{Dt}$, $\eta_{2t} = \sigma_S^{-1}\varepsilon_{St}$. This leads to a rewriting of (20) and (21) as

$$(22) \quad q_t^* = \sigma_D^{-1}q_t \\ = \left(1 - \frac{\beta\gamma}{(1 + \beta\gamma)}\right)\eta_{1t} - \frac{\beta\rho\eta_{2t}}{(1 + \beta\gamma)}$$

$$(23) \quad p_t^* = \sigma_S^{-1}p_t \\ = \frac{\gamma\rho^{-1}\eta_{1t}}{(1 + \beta\gamma)} + \frac{\eta_{2t}}{(1 + \beta\gamma)},$$

where $\rho = \frac{\sigma_S}{\sigma_D}$. Now the impulse responses to ε_{jt} have exactly the same signs as those for η_{jt} but the magnitude of the latter depends upon the ratio of the standard deviations of the cost and demand shocks (ρ), and not on their separate values. Moreover, it is clear from (22) and (23) that the impulse responses for a unit shock to η_{jt} describe the effects on q_t^* and p_t^* and not on q_t and p_t . Fundamentally, the difficulty is that $\varepsilon_{Dt} = \sigma_D\eta_{1t}$ and $\varepsilon_{St} = \sigma_S\eta_{2t}$, meaning that the η_{jt} are not equal to the demand and supply shocks, but are scaled versions of them. Another way of describing the significance of this is that the impulse responses to unit shocks in η_{jt} indicate the responses of q_t and p_t to *one standard deviation shocks* in ε_{Dt} and ε_{St} . Hence, in partial answer to the question of this subsection, correct impulse responses to one standard deviation demand and cost shocks should be recoverable using sign restrictions (provided the summative model is correct).

When would we be happy to have just one standard deviation responses, i.e., those provided by sign restriction information? Two cases come to mind. One is if we are looking at when a peak in impulses occurs, e.g., whether there is overshooting in exchange rates as in Scholl and Uhlig (2008), since the location of the peak is invariant to any positive scaling factor for the impulses. Another

would be when variance decompositions are being computed, since here what is needed are impulse responses to one standard deviation shocks.

When would we be less enthusiastic about the impulses found with sign restrictions? In many policy-related contexts, we want to answer questions such as “what would be the responses to a 100 basis point interest rate shock” or, in the market model context, to a unit shock in demand? In the latter case, we would need to know the standard deviation of the demand shocks ε_{Dt} in order to work out an answer from the sign restriction impulse responses, as these only provide the impact of one standard deviation shocks. But, because the magnitude of the standard deviation of the demand shocks is not an estimable quantity with just sign restrictions, we cannot construct impulse responses to answer questions like those just posed (unless of course $\sigma_D = 1$). In this important sense, the sign restriction approach would not recover the required impulse responses.⁸

The discussion above has centered on whether the true impulses would be in the range of models identified by sign restrictions. Leaving aside the issue that one might need to generate a very large number of these models to ensure that, we are still left with the problem of which one to select. As we have mentioned earlier, the “median” impulses are often presented. But there is nothing that says that the true impulses would coincide with the median. One feels that often the “median” impulses are thought of as “most probable,” but, as we pointed out earlier, the range of

⁸It might be that the differences between the impulse responses seen in figure 3 come from the fact that the standard errors of shocks estimated from the NK model are inaccurate due to specification problems with that model, and so these cannot be compared with the one unit shock responses in the η_{jt} found from the sign restrictions. As mentioned, the latter would be equivalent to doing one standard deviation shocks from the “true” model (provided it is described by a VAR(1)).

impulses is due to multiple models and not any uncertainty coming from data.

To emphasize that multiple models create problems in deciding on the values of the true impulses, we perform the following experiment. Suppose a structural macro model is used to generate a very large sample of data, and the model is designed so that it has impulse responses that agree with the sign restrictions in table 3. In the analysis of section 3, alternative impulses were found by recombining those for the shocks in an arbitrary recursive model. We termed these the *base shocks*. Instead, let us take the base shocks to be those from the structural macro model itself. These certainly qualify, being uncorrelated, although they would not be available in practice. Nevertheless, if we were fortunate enough to know them, the impulse responses generated by sign restrictions will be combinations of the true ones, with weights given by Q . When $Q = I$, we will get the true impulse responses. So where in the range of models do the true impulse responses lie? In our scenario, the true impulse responses were chosen to obey the table 3 sign restrictions and were close to those of the empirically estimated NK model. Computing a range of estimates by choosing different Q , it was found that the true impact impulse responses of output and inflation to an interest rate shock lay at the 12.5 and 0.4 percentiles, far from the median. Hence, the implication of this experiment is that, even if the true impulses lie in the range of models generated by the sign restrictions, we do not know where in the range they are. All we can say is that, if the range is very narrow, then we should get a good indication of what the true impulses are. Otherwise we cannot know.

Paustian (2007) performs Monte Carlo experiments on models where sign restrictions on a set of (primary) variables are imposed to identify the shocks, and then the impulses to other (secondary) variables are checked to see if they have the correct sign. He draws two conclusions from the experiments. Firstly, it

is likely that the correct signs for the impact of the shocks on the secondary variables will be found if the identified shocks have a dominant influence on the primary variables. Secondly, the more shocks that are identified the greater is the likelihood that the correct signs will be recovered. This leads him to conclude that sign restrictions can reliably recover some *qualitative* features of impulse responses under certain conditions.

The results he gets can be explained. Because the reduced form VAR shocks are e_t , and the structural ones are ε_t , the connection between them is $e_t = B_0^{-1}\varepsilon_t$. If there are no lags and $n = 3$, the first "VAR" equation will have the following relation between its error and the structural shocks:

$$(24) \quad e_{1t} = b_0^{11}\varepsilon_{1t} + b_0^{12}\varepsilon_{2t} + b_0^{13}\varepsilon_{3t},$$

where b_0^{ij} are the coefficients of B_0^{-1} . If ε_{1t} was known, then b_0^{11} (the impact response) could be consistently estimated by regressing e_{1t} on ε_{1t} , since the omitted regressors $\varepsilon_{2t}, \varepsilon_{3t}$ are uncorrelated with ε_{1t} . However, ε_{1t} is not known and sign restrictions involve combining the VAR errors e_{jt} with weights to extract an estimate ε_{it}^* . Such an estimate can be written as a combination of the ε_{jt} :

$$(25) \quad \varepsilon_{1t}^* = \phi_1\varepsilon_{1t} + \phi_2\varepsilon_{2t} + \phi_3\varepsilon_{3t}.$$

From (25) it is clear that a regression of e_{1t} on ε_{1t}^* will produce a biased estimator of b_0^{11} owing to the simultaneous presence of $\varepsilon_{2t}, \varepsilon_{3t}$ in the regressor and the error term of the equation. Of course this bias will decline as the variance of ε_{1t}^* increases relative to the variance of $b_0^{12}\varepsilon_{2t} + b_0^{13}\varepsilon_{3t}$, and this is the first conclusion Paustian reaches.

To see the second, we just need to note that, if a second shock is identified, the regression becomes one of e_{1t} on ε_{1t}^* and ε_{2t}^* . There is no certainty, but it is likely that the biases will be smaller now than before. If it was the case that ε_{2t} had been correctly estimated then it would have been eliminated from the error

term of the regression, leaving only ε_{3t} . So it is likely that, as we estimate more shocks using sign restrictions, the bias will be reduced. Again however this is not a general result as it depends upon the extent to which the shocks have been correctly extracted.

5.2 Other Summative Representations

As mentioned in the introduction to the section, if we try to align theory-inspired interpretative models (such as DSGE models) with the summative model, we often encounter the situation that there are variables in the former that are not observable, and so the latter model is fitted with a smaller number of variables.⁹ Let the observable variables (data) be z_t and the larger set in the theoretical model be z_t^+ . Then, it has been known for a long time, see Kenneth F. Wallis (1977) and Arnold Zellner and Franz Palm (1974), that a VAR in z_t^+ becomes a vector autoregressive moving average (VARMA) in z_t . Thus a VAR will not represent the data precisely if it should be generated by a theoretical model with latent (unobserved) variables, although, if one makes the order sufficiently high, it might be argued to be approximately correct. Basically this implies that the impulse responses from the theoretical model, C_j^+ , will not be equal to those from an approximating VAR, unless the order is infinite. As shown in George Kapetanios, Pagan, and Alasdair Scott (2007), this difference can be very large for some shocks and models and so one needs to exercise care in using information from theory-consistent models to identify shocks in VARs.¹⁰

⁹Technology is an obvious example of a variable in a DSGE model that is rarely present in an estimated VAR. But it is also the case that researchers often treat the capital stock as unobservable and so it is omitted from the list of variables in the empirical VAR.

¹⁰Using a model that was a smaller version of the Bank of England Quarterly Model but a VAR with only a standard set of six variables, Kapetanios, Pagan, and Scott (2007) found that a VAR(50) and thirty thousand observations were needed to recover the true impulse responses.

Of course it is possible that this problem is less of an issue for the signs of the responses than it is for the magnitudes, i.e., the signs of C_j^+ and C_j may agree even if the magnitudes do not.¹¹ Fundamentally, the problem is that a VAR is not the correct summative model.¹² As an alternative one might estimate a VARMA process or a VAR with some latent variables, but mostly researchers have dealt with the latent variable problem by expressing the theoretic model in a state space form (SSF)

$$(26) \quad z_t = H z_t^+$$

$$(27) \quad z_t^+ = M z_{t-1}^+ + G \varepsilon_t^+,$$

and then estimating this. Readily available computer programs such as Dynare are designed to do so. Thus the role of a theory-inspired model is to provide the variables in z_t^+ , and the order of the VAR associated with them, while the empirical investigator selects z_t . In the DSGE model, M and G will be functions of the model parameters θ (H simply selects variables and so generally doesn't depend on θ). The appropriate summative model therefore is the SSF (26)–(27), but with M being treated as unrestricted and $G \varepsilon_t^+$ being replaced by some errors e_t . Once estimated the residuals \hat{e}_t can be combined together using the Q matrices dealt with earlier to produce new shocks η_t and then passed through the estimated SSF to find the impulse responses for these new shocks. As noted in footnote 1, it will generally be easiest to produce initial η_t shocks that are uncorrelated by performing a Choleski decomposition upon \hat{e}_t , but an alternative approach would be to make M triangular,

¹¹Indeed this seems to be supported by the simulations in Canova and Paustian (2007), although it may just reflect the particular context they are working in.

¹²There are even cases where there is no invertible MA representation, and so no VAR exists.

and to then estimate the resulting SSF by a MLE. A program such as Dynare would enable one to do this efficiently.

5.3 Permanent and Transitory Shocks

If there are as many permanent shocks to be identified as there are observable variables, then this would imply that there is no cointegration between the variables. Therefore, the appropriate summative model is a VAR in differenced variables. Hence it is only how the data is measured that changes, allowing sign restrictions to be easily imposed by working on the residuals from the differenced-variables VAR. Sometimes one sees such a summative model in the sign restrictions literature. Examples are Marek Jarociński and Frank R. Smets (2008) and Katie Farrant and Peersman (2006). But it needs to be stressed that all shocks have to be regarded as transitory for this summative model to be correct. When there are both transitory and permanent shocks there is cointegration, and so the summative model will be the VECM

$$\Delta z_t^+ = \alpha \beta' z_{t-1}^+ + e_t.$$

Correspondingly, a structural VECM (SVECM) of the form

$$(28) \quad B_0 \Delta z_t^+ = (B_0 \alpha) \beta' z_{t-1}^+ + \varepsilon_t,$$

would be used to interpret the data. Because this is a relatively simple extension of the standard approach, it does not need any extensive development. It is only if there are latent variables that special issues arise. Hence, if z_t rather than z_t^+ is observed (and z_t has less variables than z_t^+), the problems identified in the preceding subsection arising from latent variables occur again, although they can be solved in the same way, namely via a state space form. This situation arises in many DSGE models in which the permanent

shock is technology while all other shocks are transitory.

Care does need to be exercised in finding the structural shocks. There are standard formulae for converting the VECM residuals \hat{e}_t into $n - r$ permanent \hat{e}_t^P and r transitory \hat{e}_t^T uncorrelated shocks, where r is the degree of cointegration. It would then be necessary to recombine these with Q matrices to produce new uncorrelated permanent and transitory structural shocks $\hat{\eta}_t^P$ and $\hat{\eta}_t^T$. In doing so, we need to recognize that one cannot combine the permanent shocks to produce a transitory shock. Consequently, it is simplest to work with Q_P, Q_T (each being Givens or QR) such that $\hat{\eta}_t^P = Q_P \hat{e}_t^P$ and $\hat{\eta}_t^T = Q_T \hat{e}_t^T$. To initialize the sequence, one could begin with a recursive SVAR in which $n - r$ of the structural shocks are designated to be permanent and r to be transitory. The methodology outlined in Pagan and M. Hashem Pesaran (2008) illustrates how such a system can be constructed and estimated.

6. Conclusion

When sign restriction work first began, it was mainly about the identification of a single shock. Since then it has become popular to identify multiple shocks. Moreover, the range of applications has grown from the initial focus on monetary policy. Given that sign information is rather weak, we suspect that it is best to utilize the restrictions in conjunction with parametric restrictions and that seems to be an emerging tendency as well. A number of other themes also seem to be developing. One is that contemporaneous restrictions might be preferable to imposing restrictions on longer lags. Another is that DSGE models are a useful way of finding out likely sign restrictions, particularly as the number of variables in the VAR grows. We have tried to show these tendencies in

the review, and have also argued that more care often needs to be taken in devising the model that is to summarize the data, a clear statement about whether shocks are permanent or transitory should be provided, and an account of how the multiple models and multiple shocks problems were dealt with must be present in the research. In some instances, the latter have not been well understood and often the responses to them have been not well documented.

Table 1 provides a quick summary of the studies that appear in the literature, characterized by a number of the items mentioned above—viz. whether there are a mixture of sign and other restrictions, namely whether there are permanent shocks, how many shocks are identified and whether the source of the restrictions comes from informal ideas or from a formal theory-oriented model. This provides a quick overview of the diversity of the studies. As well we classify them according to the main issue being dealt with such as the isolation of the effects of technology shocks, monetary policy shocks, and fiscal policy shocks. It is apparent from this table that sign restrictions have become of increasing interest to applied researchers seeking information about a large range of phenomena.

On balance, we do feel that sign restrictions have provided a useful technique for quantitative analysis. There are a number of instances in which variables are simultaneously determined and it is hard to justify any parametric restrictions to resolve the identification problem. A classic example is that of interest and exchange rates. In these cases, sign restrictions appeal. In other situations, such as isolating monetary policy, it seems more likely that using institutional knowledge to provide parametric restrictions would be a better way to proceed. This points to the fact that combinations of restrictions are likely to be what we will need to adopt in the future to carry out good applied work.

REFERENCES

- An, Lian. 2006. "Exchange Rate Pass-Through: Evidence Based on Vector Autoregressions with Sign Restrictions." Unpublished.
- Bjørnland, Hilde C., and Jørn I. Halvorsen. 2008. "How Does Monetary Policy Respond to Exchange Rate Movements? New International Evidence." Norges Bank Working Paper 15.
- Blanchard, Olivier Jean, and Danny Quah. 1989. "The Dynamic Effects of Aggregate Demand and Supply Disturbances." *American Economic Review*, 79(4): 655–73.
- Canova, Fabio, and Gianni De Nicoló. 2002. "Monetary Disturbances Matter for Business Fluctuations in the G-7." *Journal of Monetary Economics*, 49(6): 1131–59.
- Canova, Fabio, and Matthias Paustian. 2010. "Measurement with Some Theory: A New Approach to Evaluate Business Cycle Models." Universitat Pompeu Fabra Department of Economics and Business Working Paper 1203.
- Chari, V. V., Patrick J. Kehoe, and Ellen R. McGrattan. 2008. "Are Structural VARs with Long-Run Restrictions Useful in Developing Business Cycle Theory?" *Journal of Monetary Economics*, 55(8): 1337–52.
- Cho, Seonghoon, and Antonio Moreno. 2006. "A Small-Sample Study of the New-Keynesian Macro Model." *Journal of Money, Credit, and Banking*, 38(6): 1461–81.
- Dedola, Luca, and Stefano Neri. 2006. "What Does a Technology Shock Do? A VAR Analysis with Model-Based Sign Restrictions." European Central Bank Working Paper 705.
- Dungey, Mardi, and Renée Fry. 2009. "The Identification of Fiscal and Monetary Policy in a Structural VAR." *Economic Modelling*, 26(6): 1147–60.
- Eickmeier, Sandra, Boris Hofmann, and Andreas Worms. 2009. "Macroeconomic Fluctuations and Bank Lending: Evidence for Germany and the Euro Area." *German Economic Review*, 10(2): 193–223.
- Farrant, Katie, and Gert Peersman. 2006. "Is the Exchange Rate a Shock Absorber or a Source of Shocks? New Empirical Evidence." *Journal of Money, Credit, and Banking*, 38(4): 939–61.
- Faust, Jon. 1998. "The Robustness of Identified VAR Conclusions about Money." *Carnegie-Rochester Conference Series on Public Policy*, 49: 207–44.
- Francis, Neville R., Michael T. Owyang, and Jennifer E. Roush. 2005. "A Flexible Finite-Horizon Identification of Technology Shocks." Board of Governors of the Federal Reserve System International Finance Discussion Paper 832.
- Francis, Neville R., Michael T. Owyang, and Athena T. Theodorou. 2003. "The Use of Long-Run Restrictions for the Identification of Technology Shocks." *Federal Reserve Bank of St. Louis Review*, 85(6): 53–66.
- Fry, Renée, and Adrian Pagan. 2005. "Some Issues in Using VARs for Macroeconometric Research." Australian National University Centre for Applied

- Macroeconomic Analysis Working Paper 19.
- Fry, Renée, and Adrian Pagan. 2007. "Some Issues in Using Sign Restrictions for Identifying Structural VARs." National Centre for Econometric Research Working Paper 14.
- Fujita, Shigeru. 2011. "Dynamics of Worker Flows and Vacancies: Evidence from the Sign Restriction Approach." *Journal of Applied Econometrics*, 26(1): 89–121.
- Gali, Jordi. 1992. "How Well Does the ISLM Model Fit Postwar U.S. Data." *Quarterly Journal of Economics*, 107(2): 709–38.
- Hau, Harald, and Helene Rey. 2004. "Can Portfolio Rebalancing Explain the Dynamics of Equity Returns, Equity Flows, and Exchange Rates?" *American Economic Review*, 94(2): 126–33.
- Jääskelä, Jarkko, and David Jennings. 2010. "Monetary Policy and the Exchange Rate: Evaluation of VAR Models." Reserve Bank of Australia Economic Research Department Research Discussion Paper 2010-07.
- Jarociński, Marek, and Frank R. Smets. 2008. "House Prices and the Stance of Monetary Policy." *Federal Reserve Bank of St. Louis Review*, 90(4): 339–65.
- Kapetanios, George, Adrian Pagan, and Alasdair Scott. 2007. "Making a Match: Combining Theory and Evidence in Policy-Oriented Macroeconomic Modeling." *Journal of Econometrics*, 136(2): 565–94.
- Kilian, Lutz, and Dan Murphy. 2009. "Why Agnostic Sign Restrictions Are Not Enough: Understanding the Dynamics of Oil Market VAR Models." Unpublished.
- Leamer, Edward E. 1978. *Specification Searches and Ad Hoc Inference with Nonexperimental Data*. New York: Wiley.
- Lewis, Vivien J. 2007. "Productivity and the Euro–Dollar Real Exchange Rate." *Review of World Economics/Weltwirtschaftliches Archiv*, 143(2): 324–48.
- Mountford, Andrew. 2005. "Leaning into the Wind: A Structural VAR Investigation of UK Monetary Policy." *Oxford Bulletin of Economics and Statistics*, 67(5): 597–621.
- Mountford, Andrew, and Harald Uhlig. 2005. "What Are the Effects of Fiscal Policy Shocks?" Humboldt University Collaborative Research Center Discussion Paper 2005-039 SFB 649.
- Mountford, Andrew, and Harald Uhlig. 2008. "What Are the Effects of Fiscal Policy Shocks?" National Bureau of Economic Research Working Paper 14551.
- Pagan, Adrian, and M. Hashem Pesaran. 2008. "Econometric Analysis of Structural Systems with Permanent and Transitory Shocks." *Journal of Economic Dynamics and Control*, 32(10): 3376–95.
- Paustian, Matthias. 2007. "Assessing Sign Restrictions." *B.E. Journal of Macroeconomics*, 7(1).
- Peersman, Gert. 2005. "What Caused the Early Millennium Slowdown? Evidence Based on Vector Autoregressions." *Journal of Applied Econometrics*, 20(2): 185–207.
- Peersman, Gert, and Roland Straub. 2009. "Technology Shocks and Robust Sign Restrictions in a Euro Area SVAR." *International Economic Review*, 50(3): 727–50.
- Preston, Alan J. 1978. "Concepts of Structure and Model Identifiability for Econometric Systems." In *Stability and Inflation: A Volume of Essays to Honour the Memory of A. W. H. Phillips*, ed. A. R. Bergstrom, A. J. L. Catt, M. H. Peston, and B. D. J. Silverstone, 275–97. New York: Wiley.
- Quenouille, Maurice H. 1957. *The Analysis of Multiple Time-Series*. London: Griffin.
- Rafiq, M. Sohrab, and Sushanta K. Mallick. 2008. "The Effect of Monetary Policy on Output in EMU3: A Sign Restriction Approach." *Journal of Macroeconomics*, 30(4): 1756–91.
- Rubio-Ramirez, Juan Francisco, Daniel Waggoner, and Tao Zha. 2005. "Markov-Switching Structural Vector Autoregressions: Theory and Application." Federal Reserve Bank of Atlanta Working Paper 2005-27.
- Ruffer, Rasmus, Marcelo Sanchez, and Jian-Guang Shen. 2007. "Emerging Asia's Growth and Integration—How Autonomous Are Business Cycles?" European Central Bank Working Paper 715.
- Sanchez, Marcelo. 2007. "What Drives Business Cycles and International Trade in Emerging Market Economies?" European Central Bank Working Paper 730.
- Scholl, Almuth, and Harald Uhlig. 2008. "New Evidence on the Puzzles: Results from Agnostic Identification on Monetary Policy and Exchange Rates." *Journal of International Economics*, 76(1): 1–13.
- Sims, Christopher A. 1980. "Macroeconomics and Reality." *Econometrica*, 48(1): 1–48.
- Straub, Roland, and Gert Peersman. 2006. "Putting the New Keynesian Model to a Test." International Monetary Fund Working Paper 06/135.
- Uhlig, Harald. 2005. "What Are the Effects of Monetary Policy on Output? Results from an Agnostic Identification Procedure." *Journal of Monetary Economics*, 52(2): 381–419.
- Vargas-Silva, Carlos. 2008. "Monetary Policy and the US Housing Market: A VAR Analysis Imposing Sign Restrictions." *Journal of Macroeconomics*, 30(3): 977–90.
- Wallis, Kenneth F. 1977. "Multiple Time Series Analysis and the Final Form of Econometric Models." *Econometrica*, 45(6): 1481–97.
- Wold, Herman O. A. 1951. "Dynamic Systems of the Recursive Type—Economic and Statistical Aspects." *Sankhyā*, 11(3–4): 205–17.
- Zellner, Arnold, and Franz Palm. 1974. "Time Series Analysis and Simultaneous Equation Econometric Models." *Journal of Econometrics*, 2(1): 17–54.

This article has been cited by:

1. James H. Stock¹ and Mark W. Watson² ¹James H. Stock is the Harold Hitchings Burbank Professor of Political Economy, Harvard University, Cambridge, Massachusetts. james_stock@harvard.edu ²Mark W. Watson is Howard Harrison and Gabrielle Snyder Beck Professor of Economics and Public Affairs, Princeton University, Princeton, New Jersey. mwatson@princeton.edu . 2017. Twenty Years of Time Series Econometrics in Ten Pictures. *Journal of Economic Perspectives* 31:2, 59-86. [[Abstract](#)] [[View PDF article](#)] [[PDF with links](#)]
2. Christiane Baumeister, Gert Peersman. 2013. Time-Varying Effects of Oil Supply Shocks on the US Economy. *American Economic Journal: Macroeconomics* 5:4, 1-28. [[Abstract](#)] [[View PDF article](#)] [[PDF with links](#)]