

Review Problem Answers

1. Saying the residual at each t is drawn from a mixture of normals means that the p.d.f. of $e(t)$ at each t is the sum of three normal p.d.f.'s with weights q_i summing to one. But this can be modeled as involving the unobservable random variable i_t , which takes on the values 1, 2, 3 with probabilities q_i , $i=1, 2, 3$, and with $e(t)|i_t \sim N(0, \Sigma_{i_t})$. Collecting lagged y 's into a single row vector $X(t)$ and the coefficients in B into a single parameter matrix b , we can rewrite (1) from the problem set as

$$y(t)' = X(t)b + e(t)' . \quad [1]$$

Then the log of the joint p.d.f. of $\{i_t\}$ and $\{y(t)\}$ is

$$\sum_{t=1}^T \left\{ -\frac{1}{2} \log |\Sigma_{i_t}| - \frac{1}{2} (y(t)' - X(t)b) \Sigma_{i_t}^{-1} (y(t)' - X(t)b)' + \log q_{i_t} \right\} . \quad [2]$$

To form the likelihood in terms just of the parameters (not including $\{i_t\}$), we would have to integrate this p.d.f. over the $\{i_t\}$ sequences, which would give it an inconvenient, non-Gaussian form. It is therefore convenient not to integrate out the $\{i_t\}$'s, but to leave them in and use Monte Carlo methods that draw jointly from a p.d.f. proportional to [2] (exponentiated) in the parameters and $\{i_t\}$ jointly. We are not directly interested in the $\{i_t\}$'s, but drawing from the flat-prior posterior (the likelihood) is much easier if we draw them together with the parameters.

Conditional on a sequence $\{i_t\}_{t=1}^T$, each disturbance is normal, and this p.d.f. has a Gaussian shape in b . It is not a very nice Gaussian shape, because the occurrence of different Σ_{i_t} 's at different observations, if they are not diagonal, prevents application of any of the usual SUR conditions that allow equation-by-equation estimation of the peak of the likelihood. Nonetheless, it is a big Gaussian distribution. Its shape can be found by stacking up the data for the system, treating it as a big GLS problem with nT observations, where n is the number of variables in y and T is the number of observations. This will give a mean and covariance matrix for the elements of B , and it is then straightforward to draw a random B from this distribution.

Conditional on all the other parameters and the data, the p.d.f. as a function of the i_t 's is straightforward. The $e(t)$ sequence is a known function of the conditioning variables. A particular i_t enters the log likelihood only through a single element in the sum in [2], the term indexed by t . This means that the p.d.f. as a function of $\{i_t\}_{t=1}^T$ has the form of independent draws for the elements of the sequence. There are three possible values of this term, corresponding to the three possible values of i_t , so to draw from this conditional distribution for i_t requires just a simple draw from a three-point discrete distribution at each t .

Conditional on all the other parameters, the data, and the $\{i_t\}$ sequence, the three covariance matrices Σ_i , $i=1, 2, 3$ enter the joint p.d.f. through three separate components of the sum in [2]. For Σ_1 , for example, the corresponding component of [2] is

$$-\frac{T_1}{2} \log |\Sigma_1| - \frac{1}{2} \sum_{t \in I_1} (y(t)' - X(t)b) \Sigma_1^{-1} (y(t)' - X(t)b)' - T_1 \log q_1 \quad [3]$$

where T_1 is the number of observations in I_1 , the set of t values for which $i_t = 1$. Note that as a function of Σ_1 this is in the form of a standard inverse-Wishart p.d.f. One forms the residual covariance matrix for each of the three subsamples I_j , then draws a value for the corresponding Σ_j from the associated inverse-Wishart. (See Gelman, Carlin, Stern and Rubin, p.474-5 and 481, for description of the inverse-Wishart).

Looking again at [3] we can see that the p.d.f. as a function of the q 's is proportional simply to

$$q_1^{T_1} q_2^{T_2} (1 - q_1 - q_2)^{T_3} . \quad [4]$$

This is a three-dimensional Dirichlet p.d.f. (GCSR p. 476-7, 481), which is straightforward to draw from.

Thus we have the possibility of an exact Gibbs sampling scheme, drawing in turn from the conditional posteriors for b , $\{i_t\}_{t=1}^T$, $\{\Sigma_i\}_{i=1}^3$, $\{q_i\}_{i=1}^3$.

All of this discussion has assumed we would use the conditional p.d.f. of the y 's given initial observations. If we assume stationarity, there would be an implied marginal distribution for the initial y 's we could use. However, this marginal distribution is extremely messy analytically. Since Markov-chain Monte Carlo methods assume that the p.d.f. of the data can be calculated for given parameter values at given data points, they will not help with this situation.

2. Asymptotically, for the stationary case ($|r| < 1$), the usual OLS formulas for mean and variance of the OLS estimators and the usual Gaussian distribution for the estimators are justified. So when the model is true, we can write asymptotically

$$\sqrt{T}(\hat{r}_T - r) \sim N(0, 1 - r^2) . \quad [5]$$

Also, for the stationary case the error in \hat{r}_T is negligibly correlated with the last observation in the sample, so the distribution of $\hat{r}_T^2 y(T)$ conditional on $y(T)$ is asymptotically just the marginal asymptotic distribution of \hat{r}_T^2 scaled by $y(T)$. Note that the mean-squared error of two-step-ahead forecast in this "true model" case is dominated by the variance of the forecast based on knowledge of the true r , which is $(1 + r^2) S_e^2$. The component due to error in estimation of r declines at the rate $T^{-\frac{1}{2}}$. The asymptotic distribution of \hat{r}_T^2 is centered at r^2 , with variance (after scaling by \sqrt{T})

$$(2r)^2 \cdot \text{var}(\sqrt{T}(\hat{r}_T - r)) = 4r^2 \cdot (1 - r^2). \quad [6]$$

If we estimate using OLS on (3) from the problem statement, the error in the coefficient estimate is

$$\frac{\sum (e(t) + re(t-1))y(t-2)}{\sum y(t-1)^2}. \quad [7]$$

The denominator of this expression, when normalized by T^{-1} , converges in probability to the constant $s^2/(1-r^2)$, while the numerator, being a sum of zero-mean, finite-variance, stationary random variables, converges, when normalized by $T^{-\frac{1}{2}}$, to a zero-mean normal distribution. There is a standard result, which I realized as I prepared this answer sheet might not have been covered in the prerequisites for this course, that time averages of a zero-mean, finite-variance stationary process x with autocovariance function $\text{cov}(x(t), x(t-s)) = R_x(s)$ converge in distribution, when normalized by $T^{-\frac{1}{2}}$, to $N\left(0, \sum_r R_x(t)\right)$, assuming that R_x is absolutely summable. The numerator of [7] has ACF $R(t) = 0$ for $t > 1$, since the e 's are zero-mean and i.i.d. It has variance $R(0) = s_e^4 \cdot (1+r^2)/(1-r^2)$ and $R(1) = s_e^4 r^2/(1-r^2)$. Therefore the ratio in [7], normalized by \sqrt{T} , converges to a $N(0, 1+3r^2)$. (To see why, be sure you remember that one sums R over both positive and negative values of its argument.) It is not hard to check then that the asymptotic variance of this estimator, $1+3r^2$, must exceed that of the previous estimator based on a one-step-ahead regression, as given in [6]. In fact, the difference always exceeds $15/16$, for any value of r . But remember that the unnormalized difference is shrinking at rate $T^{-\frac{1}{2}}$ and becomes a negligible part of the overall MSE of forecast as $T \rightarrow \infty$.

For the other case, where the AR model is false, using the one-step-ahead OLS estimate produces an inconsistent estimator of the coefficient in the best linear predictor of $y(T+2)$ based on $y(T)$. Thus the difference between the two estimators does not dwindle as $T \rightarrow \infty$. This is the main point of this question, and had been discussed in class. On an exam, recognizing this point and giving a quick argument as to why it's true in this case would have given you most of the credit for the question. But here we proceed with the details.

If y is drawn from a first-order MA of the form

$$y(t) = n(t) + an(t-1), \quad [8]$$

then the best linear predictor of $y(t)$ based on $y(t-2)$ is just zero, since y is uncorrelated with values of itself lagged by more than one period, and its MSE is $s_e^2 \cdot (1+a^2)$. The best one-step ahead linear predictor for $y(t)$ based on $y(t-1)$ alone is $qy(t-1)$, with

$$q = R_y(1)/R_y(0) = a/(1+a^2). \quad [9]$$

If we estimate a regression of $y(t)$ on $y(t-1)$ by OLS, we will obtain a consistent estimate of q , and the resulting two-step ahead prediction error will be, in large samples

$$s_e^2 \cdot (1+a^2) + q^2 s_y^2 = s_e^2 \frac{1+3a^2+a^4}{1+a^2}, \quad [10]$$

which is clearly larger than what is obtained with the OLS regression on $y(t-2)$, except when $a = 0$. The ratio of MSE's is worst when a^2 is at its upper limit of 1, when the 1-step-ahead regression gives an MSE 1.25 times as large as the 2-step-ahead regression.

3. i) The money neutrality proposition, in a log-linear system, is a single system-wide dummy observation. It will be a scalar multiple of an observation in which all nominal variables are set to one and all real variables are set to zero. When multiplied by the coefficients, this artificial data vector implies exactly the neutrality constraint. Since the likelihood itself has the desired Kronecker structure, adding this pseudo-observation will preserve that structure. Note that we will have to accept that the variance of the disturbance is the same in all equations for this dummy observation, so we cannot be more sure of the neutrality restriction in some equations than in others.

ii) A restriction that the sum of coefficients on variable j in equation i be zero is not in the form of a system-wide dummy observation, as it involves only coefficients in one equation. However, as the Sims-Zha article points out, differences across equations in the right-hand-side variable matrix, while destroying the strict Kronecker structure, do still allow equation-by-equation estimation when the model is formulated to imply a diagonal covariance matrix of disturbances. The requirement is just that block diagonality of the covariance matrix in the posterior for a^+ be preserved. Since the long-run restriction we are discussing here is a linear combination of coefficients in a single equation, it can be expressed by adding a dummy observation to the data matrix for that equation, and this will preserve block-diagonality.

Note that most existing empirical work using “long run restrictions” on VAR's has used 2-variable models, where the restriction that an off-diagonal element of the AR operator have coefficients summing to zero is equivalent to the restriction that the corresponding element of the MA operator have coefficients summing to zero. But in systems of dimension larger than 2 this correspondence no longer holds. Often a long-run restriction is more reasonably formulated as, say, “monetary policy shocks have no long run effects on Y ” rather than as “a permanent change in M will have no long run effect on Y ”. The latter is a restriction on the AR, the former a restriction on the MA. Linear restrictions on the MA are in general complicated nonlinear restrictions on the AR, involving all equations even when the MA restriction involves only one i and j .