

Multivariate STAR Unemployment Rate Forecasts

Costas Milas*

Department of Economics, City University

Philip Rothman†

Department of Economics, East Carolina University

September 21, 2004

Abstract

In this paper we use smooth transition vector error-correction models (STVECMs) in a simulated out-of-sample forecasting experiment for the unemployment rates of the four non-Euro G-7 countries, the U.S., U.K., Canada, and Japan. For the U.S., pooled forecasts constructed by taking the median value across the point forecasts generated by the STVECMs perform better than the linear VECM benchmark more so during business cycle expansions. Pooling across the linear and nonlinear forecasts tends to lead to statistically significant forecast improvement for business cycle expansions for Canada, while the opposite is the case for the U.K.

Keywords

nonlinear, asymmetric, STVECM, pooled forecasts, Diebold-Mariano

Word Count: 9,557 words

*Department of Economics, City University, Northampton Square, London EC1V OHB UK, **email:** c.milas@city.ac.uk

†Department of Economics, Brewster Building, East Carolina University, Greenville, NC 27858 USA, **email:** rothmanp@mail.ecu.edu. Support from an East Carolina University Faculty Senate research grant is gratefully acknowledged.

1 Introduction

Applying the statistical theory of finite-state Markov chains, Neftci (1984) reported evidence showing that the U.S. quarterly unemployment rate is asymmetric in the sense that the probability of a decrease in the series, conditional on two preceding decreases, is greater than the corresponding probability of an increase conditional on two previous increases. One of the primary time series implications of such behavior is that it is inconsistent with a linear data generating process with symmetrically distributed innovations.

Neftci's study inaugurated a near two-decade long research program in which the extent to which key business cycle indicators display varying forms of asymmetric dynamics has been explored; see Clements and Krolzig (2003) for a useful survey of important developments in the business cycle asymmetry literature. Many of these papers have focused specifically on unemployment rates and have frequently documented strong evidence in favor of dynamic asymmetries, often in the form of parametric nonlinear models, for these series. Recent work includes: Altissimo and Violante (2001), who, by way of a threshold vector autoregressive (VAR) model of U.S. output and unemployment with feedback from a Beaudry and Koop (1993)-like "depth of recession" measure, identified nonlinearities in the propagation and persistence of shocks as well as a beneficial long-run effect of recessions on growth; Caner and Hansen (2001), who uncovered threshold autoregressive (TAR) effects in the U.S. unemployment rate via their TAR-based unit root test; and Skalin and Teräsvirta (2002), whose results suggest that smooth transition autoregressive (STAR) models can capture well the asymmetry displayed in many OECD unemployment rate series.

While the vast majority of these papers have been concerned with in-sample fits of linear and nonlinear models to unemployment rate data, researchers have increasingly investigated one of the main practical problems which has motivated the literature, that is, whether out-of-sample unemployment rate forecasts generated by nonlinear time series models can dominate those produced with standard linear models. Rothman (1998) was one of the first to consider this question for the U.S. quarterly unemployment rate. He analyzed the forecasting performance of six nonlinear time series models against linear

forecasts, and in many cases the mean squared prediction error (MSPE) associated with the nonlinear forecasts was less than those for the linear forecasts. Montgomery, Zarnowitz, Tsay, and Tiao (1998) found that, at multistep-ahead forecast horizons during business cycle contractions, TAR and Markov-switching autoregressive models outperformed in the MSPE-sense the benchmark linear model in out-of-sample forecasting of the U.S. quarterly unemployment rate. Using artificial neural network (ANN) and logistic STAR (LSTAR) models for a very large data set of U.S. macroeconomic time series, including the monthly unemployment rate, Stock and Watson (1999) showed that linear forecasts generally dominated the nonlinear forecasts. However, following a similar approach with an analogous data set for the Euro area, Marcellino (2002) reported much more favorable results for ANN and LSTAR forecasts; for the Euro area unemployment rates in particular, the ANN and LSTAR forecasts had lower MSPEs two and half times more often than did the linear forecasts.¹

A common feature of the nonlinear forecasts evaluated in these four papers is that they were all univariate.² This marks a significant point of departure for our paper: while we also examine nonlinear forecasts of unemployment rates, the models we use are multivariate. The macroeconomic theoretical motivation behind a multivariate approach is straightforward; through standard arguments it is reasonable to assume that the unemployment rate is interrelated with other important variables. The degree to which a particular nonlinear parameterization of these relationships can be exploited to yield improved forecast improvement is the empirical issue addressed in this paper.

To investigate this question for unemployment rates, we employ multivariate STAR models in which we impose cointegrating restrictions. In doing so, we build upon Skalin and Teräsvirta (2002), who noted that their univariate in-sample analysis can be interpreted as a first step in the specification of a multivariate STAR model of unemployment rates.

We also follow Rothman, van Dijk, and Franses (2001), who used a similar approach to study the Granger-causal relationship between money and output. These authors found strong evidence in favor of STAR-type nonlinearity in a system of output, prices, interest rates, and money. By Okun's Law, comparable results arguably are expected to hold for an analogous model in which output is replaced by the unemployment rate. In addition to

our primary focus on unemployment rates, there are several differences between our paper and Rothman *et al.* (2001).

First, our chief concern is evaluation of the out-of-sample forecasting performance of the models, while Rothman *et al.* (2001) concentrated on both in-sample and out-of-sample results to analyze the money-output relationship. Our main in-sample interest is identification of the *transition variables* which govern parameter variation in STAR models. Second, ours is a closer approximation to real-time implementation of these forecasting models. In Rothman *et al.* (2001) specification of the STAR models was done using practically the full sample, such that common specifications were imposed in all rolling windows of data. While this aided interpretation of the results with respect to the Granger causality question under consideration, it effectively allowed the use of post-sample information in generating the forecasts. In contrast, we specify the models for each data window only using data available through the date of each forecast, and thus allow the model specifications to vary across data windows. Though this substantially increases the computational burden, we feel our experimental design offers a better simulation of real-time forecasting practice.³ Third, in this paper forecasts are computed using two approaches: following the standard route by iterating forward estimated one-step-ahead models; and also, following Stock and Watson (1999) and Marcellino (2002), by estimating directly h -step-ahead models and projecting them forward. This allows a useful comparison of these strategies for forecasting unemployment rates. Rothman *et al.* (2001) did not employ “ h -step-ahead projections” for multistep-ahead forecasting. Fourth, we consider some easily-constructed pooled forecasts, whereas Rothman *et al.* (2001) did not use any forecast pooling procedures. Finally, while Rothman *et al.* (2001) only worked with U.S. data, we examine multivariate STAR models with data for the U.S., U.K., Canada, and Japan.

Another paper quite close to ours is Krolzig, Marcellino, and Mizon (2002), who analyzed a Markov-switching vector error correction model (MS-VECM) of the U.K. labor market with quarterly data. Besides our use of STAR as opposed to MS models, there are several differences between our paper and Krolzig *et al.* (2002). First, their unemployment measure was the volume of unemployment, whereas we use unemployment rates. Second, selection of the variables to be included in their system came out of specific focus on the

labor market; their four-variable system comprised unemployment, employment, real output, and real wages. In contrast, our choice of variables follows a standard practice in the empirical monetary policy literature; our four-variable system comprises the unemployment rate, the aggregate price level, a monetary aggregate, and a short-term interest rate. Third, in their out-of-sample forecasting exercise, Krolzig *et al.* (2002) only computed one-step-ahead forecasts, and do so only for two estimated versions of their model.⁴ In our approach the models used are reestimated for each fixed-length rolling window of data, and we compute one-quarter-ahead through eight-quarters-ahead forecasts.

The paper proceeds as follows. In Section 2 we discuss the multivariate STAR model and outline a specification procedure for such models. The results of linearity testing against STAR alternatives within a multivariate context are also presented in this section. Our out-of-sample forecasting results are examined in Section 3 and Section 4 concludes the paper.

2 Multivariate STAR Models and Linearity Testing

Let $\mathbf{x}_t = (x_{1t}, \dots, x_{kt})'$ be a $(k \times 1)$ vector time series. In our case we have $\mathbf{x}_t = (u_t, m_t, p_t, i_t)'$, with u_t the log of the unemployment rate, m_t the log of a money supply measure, p_t the log of the producer price index, and i_t a short-term interest rate.⁵ We analyze quarterly vector time series for four different countries, the U.S., U.K., Canada, and Japan for the 1959:1-2001:4, 1965:1-2001:4, 1968:1-2001:4, and 1966:4-2001:4 sample periods, respectively.⁶ The unemployment rate, money supply, and producer price index series were seasonally adjusted, while the interest rate series were not. The data were obtained from the following sources: the Federal Reserve Bank of St. Louis, the U.S. Bureau of Labor Statistics, the U.K. Office for National Statistics, the Bank of Canada, and the OECD *Main Economic Indicators* and the IMF *International Financial Statistics* databases.

A k -dimensional smooth transition vector error-correction model [STVECM] can be

specified as

$$\begin{aligned} \Delta \mathbf{x}_t = & \left(\boldsymbol{\mu}_1 + \boldsymbol{\alpha}_1 \mathbf{z}_{t-1} + \sum_{j=1}^{p-1} \boldsymbol{\Phi}_{1,j} \Delta \mathbf{x}_{t-j} \right) (1 - G(s_t; \gamma, c)) \\ & + \left(\boldsymbol{\mu}_2 + \boldsymbol{\alpha}_2 \mathbf{z}_{t-1} + \sum_{j=1}^{p-1} \boldsymbol{\Phi}_{2,j} \Delta \mathbf{x}_{t-j} \right) G(s_t; \gamma, c) + \boldsymbol{\varepsilon}_t, \quad (1) \end{aligned}$$

where Δ_j denotes the j -th difference operator, defined as $\Delta_j x_t = x_t - x_{t-j}$ for integers $j \neq 0$ and $\Delta_1 \equiv \Delta$, $\boldsymbol{\mu}_i$, $i = 1, 2$, are $(k \times 1)$ vectors, $\boldsymbol{\alpha}_i$, $i = 1, 2$, are $(k \times r)$ matrices, $\mathbf{z}_t = \boldsymbol{\beta}' \mathbf{x}_t$ for some $(k \times r)$ matrix $\boldsymbol{\beta}$ denoting the error-correction terms, $\boldsymbol{\Phi}_{i,j}$, $i = 1, 2$, $j = 1, \dots, p-1$, are $(k \times k)$ matrices, and $\boldsymbol{\varepsilon}_t = (\varepsilon_{1t}, \dots, \varepsilon_{kt})$ is a k -dimensional vector white noise process with mean zero and $(k \times k)$ covariance matrix $\boldsymbol{\Sigma}$. The transition function $G(s_t; \gamma, c)$ is assumed to be a continuous function bounded between zero and one. In this paper we allow the transition variable s_t to be either a function of a lagged component of \mathbf{x}_t or a lagged exogenous variable.

The STVECM can be thought of as a regime-switching model that allows for two regimes associated with the extreme values of the transition function, $G(s_t; \gamma, c) = 0$ and $G(s_t; \gamma, c) = 1$, where the transition from one regime to the other is smooth. In this paper we restrict attention to the logistic transition function

$$G(s_t; \gamma, c) = \frac{1}{1 + \exp\{-\gamma(s_t - c)/\hat{\sigma}_s\}}, \quad \gamma > 0, \quad (2)$$

where $\hat{\sigma}_s$ is the sample standard deviation of s_t . The parameter c in (2) can be interpreted as the threshold or border between the two regimes, in the sense that the logistic function changes monotonically from 0 to 1 as s_t increases, and $G(c; \gamma, c) = 0.5$. The parameter γ determines the smoothness of the change in the value of the logistic function and, thus, the smoothness of the transition from one regime to the other. As γ becomes very large, the logistic function approaches the indicator function $I[s_t > c]$. Hence, the STVECM (1) with (2) nests a two-regime threshold vector error-correction model [TVECM] as a special case; see Balke and Fomby (1997) and Tsay (1998) for discussion. Finally, note that when $\gamma = 0$ the logistic function equals 0.5 for all s_t , such that the STVECM model reduces to

a linear VECM.

The procedure we follow for specifying STVECMs is a straightforward modification of the specification procedure for univariate STAR models put forward by Teräsvirta (1994). We start by specifying a linear VECM for \mathbf{x}_t , that is,

$$\Delta \mathbf{x}_t = \boldsymbol{\mu} + \boldsymbol{\alpha} \mathbf{z}_{t-1} + \sum_{j=1}^{p-1} \boldsymbol{\Phi}_j \Delta \mathbf{x}_{t-j} + \boldsymbol{\varepsilon}_t, \quad (3)$$

where the lag order p should be such that the residuals $\hat{\boldsymbol{\varepsilon}}_t$ are approximately white noise and have zero autocorrelations at all lags. To reduce the number of parameters ($4 + (4 \times r) + (4 \times 4 \times (p - 1))$), we decided to use a subset VECM by imposing zero restrictions on coefficients in the $\boldsymbol{\Phi}_j$, $j = 1, \dots, p - 1$, matrices in (3). Use of such subset models simplifies computation of the test statistics required for the linearity tests described below and significantly eases estimation of the STVECMs used.

The subset VECM is specified by following the strategy recommended by Brüggemann and Lütkepohl (2000), which treats the individual equations in the VECM separately. We estimate the parameters in the i -th equation of (3) by ordinary least squares [OLS] and sequentially delete the regressor with the smallest absolute value of the corresponding t -ratios, until all t -ratios of the remaining coefficients are greater than some threshold value τ in absolute value; in each iteration only a single regressor is eliminated. Then the reduced model equation is re-estimated and new t -ratios are computed. We choose the threshold τ as a function of the iteration l as

$$\tau = \tau_l = \sqrt{(\exp(\lambda_T/T) - 1)(T - L + l - 1)}, \quad (4)$$

where T denotes the effective sample size, $L = 1 + r + 4 \times (p - 1)$ is the number of parameters in the unrestricted equation and λ_T is a sequence indexed by the sample size. As shown by Brüggemann and Lütkepohl (2000), by setting λ_T equal to the penalty term involved in an information criterion of choice, this procedure leads to the same final model as sequentially removing those regressors whose elimination yields the largest improvement in the value of this particular information criterion. We use the Akaike Information Criterion (AIC),

which requires setting $\lambda_T = 2$.⁷ It should be noted that we only eliminate lagged first differences from the VECM, and always retain the intercept and error-correction terms.

We set the cointegrating rank $r = 2$ and pre-specify the two cointegrating vectors as $(1, 0, 0, 0)'$ and $(0, 0, 0, 1)'$, that is, the first row of \mathbf{z}_t is the log-unemployment rate and the second row of \mathbf{z}_t is the short-term interest rate. Such pre-specification as opposed to estimation of the cointegrating vectors serves as a simplifying pair of assumptions which allows us to focus on the value-added of allowing STAR-type effects in a multivariate forecasting model of the unemployment rate. In using u_t as an error-correction term, we follow Skalin and Teräsvirta (2002), who assumed that the unemployment rate, a bounded variable, is a globally stationary process. By way of an LSTAR specification, however, asymmetry and local nonstationarity are possible.⁸ Our decision to include i_t as an error-correction term follows Rothman *et al.* (2001), who did so in the ‘Hendry-style’ in that it was based on economic theory; see, for example, Hendry and Mizon (1993) and Söderlind and Vredin (1996). The latter authors showed that the Cooley and Hansen (1995) monetary equilibrium business cycle model implies that the nominal interest rate is stationary.

The next step in the specification procedure is to select a transition variable s_t , which is done via linearity testing of the subset VECM against the alternative of a STVECM. To carry out our forecasting exercise we require a sequence of transition variables for “rolling” fixed-length windows of data, where the first data window runs from the first observation of the data set for each country out to 1991:4, and each successive data window is constructed by shifting the preceding window ahead by one observation. This setup allows us to generate 33 out-of-sample forecasts at forecast horizons $h = 1, \dots, 8$.

Testing linearity in this context is complicated by the fact that the STVECM contains nuisance parameters which are not identified under the null hypothesis; see, for example, Davies (1987). To circumvent this identification problem, we follow the approach of Luukkonen, Saikkonen, Teräsvirta (1988) and replace the transition function $G(s_t; \gamma, c)$ with a suitable Taylor approximation. The S_1 test is a standard variable addition test based on an auxiliary regression of the residuals from the linear VECM on a set of variables given by a first-order Taylor expansion of $G(s_t; \gamma, c)$. The S_2 test is based on a

third-order Taylor approximation of the logistic transition function, and the S_3 test is a parsimonious version of the S_2 test. For the sample sizes we have, it turns out that we lack sufficient degrees of freedom to compute the system-wide version of the S_2 test, such that we only use its parsimonious version.⁹

Given that our VECM residuals tend to be highly heteroskedastic, it makes sense to employ heteroskedasticity-robust versions of the linearity tests; simulations discussed by Rothman *et al.* (2001) showed that the estimated sizes of the non-robust linearity tests in the presence of heteroskedasticity tend to be severely distorted upwards. To this effect, the specification tests developed by Wooldridge (1991) are very helpful, since they can be used in the presence of heteroskedasticity without the need to specify the often unknown form of heteroskedasticity explicitly. The robust versions of the linearity tests we use were obtained by applying ‘Procedure 3.1’ of Wooldridge (1991).

Simulations discussed by Lundbergh and Teräsvirta (1998) and Rothman *et al.* (2001) suggest, however, that these robust tests are conservative, with estimated sizes less than nominal significance levels and low estimated power. Nonetheless, we follow Rothman *et al.* (2001) and apply these heteroskedasticity-robust versions of the linearity tests since we feel that the ranking across a set of prospective transition variables is valuable information for the STVECM specification process. It is unlikely that such a ranking will be affected by presence of heteroskedasticity in the VECM residuals.

To identify an appropriate transition variable s_t with a linearity test for each data window, we run the test for several candidates, s_{1t}, \dots, s_{mt} , and select the one for which the p -value of the associated test statistic is smallest. Here we consider the following different candidate transition variables for all countries: lagged yearly changes in the log unemployment rate ($\Delta_4 u_{t-d}$), lagged yearly growth rates in the money supply ($\Delta_4 m_{t-d}$), lagged annual inflation rates ($\Delta_4 p_{t-d}$), lagged yearly changes in the short-term interest rate ($\Delta_4 i_{t-d}$), lagged yearly changes in the annual money supply growth rates ($\Delta_4 \Delta_4 m_{t-d}$), lagged yearly changes in the annual inflation rate ($\Delta_4 \Delta_4 p_{t-d}$), and lagged yearly changes in the relative price of oil ($\Delta_4 o_{t-d}$, with $o_t = p_t^{\text{OIL}}/p_t$ and p_t^{OIL} the crude petroleum producer price index). In addition, for the U.S. we also used lagged annual changes in the federal funds rate ($\Delta_4 ff_{t-d}$).¹⁰

The reason why we use 4-quarter differences as transition variables is that we expect the regimes in unemployment rate dynamics to be more so persistent, because, for example, they might be related to the business cycle or to monetary policy. Using 4-quarter differences effectively eliminates short-run fluctuations which do not necessarily represent changes in regimes. We test linearity with the above-mentioned variables for delays $d = 1, \dots, d_{\max}$, where we set the maximum value of the delay parameter d_{\max} equal to 4.

The empirical and theoretical literature upon which we base our focus on these particular candidate transition variables is large. Of particular relevance in our STAR context, we note that a good deal of research has been done which suggests that these variables are reasonable measures of either the ‘state of the economy’ and/or the ‘state of policy.’ As such, our use of these variables is strongly motivated by much of the macroeconomic research on ‘state-dependent’ dynamics; see, for example, Caplin and Leahy (1991) and Caballero and Hammour (1994).¹¹

3 Out-of-Sample Forecasting

3.1 Forecasting Methods

Our forecasts are produced by 11 forecasting methods for our 33 simulated out-of-sample periods, where we use the term “methods” in the sense of Stock and Watson (1999). That is, the sequence of forecasts generated by each method is based on an underlying “primitive model,” and we let the specification of each primitive model vary across the 33 simulated in-sample periods. The first forecasting method is based on identifying a linear VECM for each rolling in-sample window, using the AIC and a diagnostic check for residual serial correlation. For each sample window the maximum lag length allowed is 4 and the model is estimated by seemingly unrelated regressions estimation. Multistep-ahead forecasts are computed by iterating forward the estimated one-step-ahead model; we end up with forecasts for steps $h = 1, \dots, 8$.

Forecasting Methods 2 through 6 use as the primitive model a STVECM, and generate multistep-ahead forecasts by, as we do with forecasting Method 1, iterating forward the estimated one-step model. But since the models are nonlinear, we use bootstrap simula-

tions to help compute the multistep-ahead forecasts. These STVECM forecasting methods differ as to how the transition variable is selected. Methods 2 through 4 use the top-ranked candidate variable as determined by the single-equation S_1 , S_2 , and S_3 tests, respectively, run on the first-differenced log-unemployment rate equation of a subset linear VECM obtained through the Brüggemann and Lütkepohl (2000) procedure. The Brüggemann and Lütkepohl (2000) algorithm is further applied to the STVECM to facilitate estimation of the model. Methods 5 and 6 do the same with the system-wide S_1 and S_3 tests, respectively.

Forecasting Methods 7 through 9 use STVECM “ h –step-ahead projections” constructed as follows. First, for each in-sample window, we estimate directly the h –step-ahead model for the log-unemployment equation of the STVECM, such that with the dependent variable Δu_t , the first lag allowed amongst the regressors is from observation $t - h$ for forecast step h . This requires that we select a transition variable for each separate forecast step h for each of the 33 in-sample rolling windows of data, which leads to selection of 792 (8 forecast steps \times 3 forecasting methods \times 33 data windows) transition variables per country. Methods 7 through 9 base the selection of the transition variable on the single-equation S_1 , S_2 , and S_3 tests, respectively, run on the first-differenced log-unemployment rate equation of an unrestricted h –step-ahead linear VECM; after selection of the transition variable, the Brüggemann and Lütkepohl (2000) procedure is applied to the STVECM log-unemployment rate equation to help identify the appropriate regressors. Second, the forecast of Δu_{t+h} is computed by projecting the estimated equation ahead by h periods.

Stock and Watson (1999) and Marcellino (2002) emphasize that use of such h –step-ahead-projections simplifies significantly computation of the multistep-ahead nonlinear forecasts, since no simulations are required for forecast steps $h > 1$. On the other hand, this requires a very large increase in the number of linearity tests run to rank the candidate transition variables for all data windows. Further, these authors point out that h –step-ahead projections can reduce the effects of misspecification of the estimated one-step-ahead, since the effects of such misspecification do not propagate through to the multistep-ahead forecasts. Estimation of all STVECMs used in Methods 2 through 9 is done by nonlinear generalized least squares.

In addition, we employ two straightforward pooling procedures. First, Method 10

forecasts are constructed by taking the median forecast value from the nonlinear forecasts produced by Methods 2 through 9. Second, Method 11 uses the median forecast across Methods 1 through 9. Table 1 summarizes all of the forecasting methods used.

3.2 Out-of-Sample Forecasting Ranks

Table 2 presents out-of-sample forecasting rankings of these methods for each of the four countries according to two evaluation criteria, the mean squared prediction error (MSPE) and median squared prediction error (MedSPE); note that the “better” or “higher ranked” forecasting methods have “lower” numerical ranks. The key result for the U.S. is that Method 10, the pooled median forecast across the STVECMs, is the top-ranked forecasting methods according to both the MSPE and MedSPE. So, in addition to dominating the linear VECM-based Method 1, median-pooling across the nonlinear forecast methods is superior to such pooling when the linear forecasts are also used. The result that median-pooling across all the nonlinear forecasting methods dominates each of the individual ones suggests that focus on single-primitive-model-based nonlinear forecasting methods may mask the potential gains obtainable by combining these individual nonlinear forecasts. Method 1, based on forecasts from the linear VECM, is the seventh-ranked forecasting method according to the MSPE criterion, and its relative performance decreases substantially, down to eleventh out of the 11 methods, using the MedSPE, the more robust forecast comparison criterion. It appears that this rather weak performance of Method 1 accounts for Method 11, which produces forecasts by taking the median point forecast across the linear and nonlinear models, being ranked fourth and eighth, respectively, according to the MSPE and MedSPE.

We next discuss the forecasting ranks of Methods 2 through 9 for the U.S., since we are interested in determining whether any particular class of STVECM forecasts used tend to dominate another. As per the definitional scheme given in Table 1, we distinguish three such classes of STVECM forecasts within Methods 2 through 9: Methods 2 through 4; Methods 5 and 6; and Methods 7 through 9. First, we note that the h -step-ahead projections of Methods 7 through 9 are outperformed, according to the MSPE, by Methods 2 through 6. However, this result does not carry through to forecast evaluation using the

MedSPE. Second, no individual non-pooling nonlinear forecasting method is dominant across use of both the MSPE and MedSPE; for example, Method 4 is the top-ranked forecasting method out of Methods 2 through 9 according to the MSPE, while Method 7 is top-ranked according to the MedSPE.

In examining the average rank results in this table, it is useful to note that if the average rank of Method i is higher than the average rank of Method j according to either the MSPE or MedSPE, then Method i outperforms Method j via the particular criterion for more than 50% of the forecast horizons, that is, for at least 5 out of the 8 forecast horizons used. We have tabulated more specific details on such pair-wise forecast method comparisons, but do not report them here in order to save space.¹²

For the U.K., Method 1 is the top-ranked forecasting method according to the MSPE and Method 10 is top-ranked using the MedSPE. So, as in the U.S. case, the relative performance of the linear VECM forecasting method worsens when the robust MedSPE evaluation criterion is used. Also, Method 10 once again dominates Method 11 according to both the MSPE and MedSPE, i.e., median-pooling is less helpful when the linear VECM forecasts are used. Using the MSPE, Method 10 is second-ranked. But using the MedSPE, Method 2, which is based on STVECM iterative multistep-ahead forecasts with the top-ranked transition variable selected by the single-equation S_1 test, is second-ranked.

Further, according to the MSPE, there is no clear ranking of the three classes of forecasts among Methods 2 through 9; Methods 2 and 4 are top-ranked, but the third member of this class, Method 3, is ranked last. That said, via the MSPE Methods 5 and 6, which select transition variables via the system-wide linearity tests, outperform the h -step-ahead projections of Methods 7 through 9. On the other hand, using the MedSPE the class of h -step-ahead projections clearly performs worst.

For Canada, Method 11 is the top-ranked forecasting method using both the MSPE and MedSPE, showing that the nonlinear forecasts provide useful information which is not incorporated in the linear VECM forecasts. The second-ranked forecasting methods are Methods 10 and 3 according to, respectively, the MSPE and MedSPE. As in the U.S. and U.K. cases, use of the MedSPE leads to a decrease in the relative performance of the linear VECM forecasting method. In addition, none of the three classes of nonlinear fore-

casts dominates another with either the MSPE and MedSPE in predicting the Canadian unemployment rate.

For Japan, the linear VECM-based Method 1 is the top-ranked forecasting method using both the MSPE and the MedSPE. So, in contrast to the results for the other three countries, for Japan the relative performance of Method 1 is constant across use of both the MSPE and MedSPE. The second-ranked forecasting method according to the MSPE is Method 11, while Method 10 is second-ranked via the MedSPE. As is the case for Canada, none of the three classes of nonlinear forecasting methods dominates another with either forecast criterion.

3.3 Statistical Significance of MSPE Reductions

3.3.1 Uniform Weighting

To examine whether the MSPE reductions we observe in Table 2 are statistically significant, we apply the Harvey, Leybourne, and Newbold (1997) modification of the Diebold and Mariano (1995) statistic (DM). The DM test statistic is computed by weighting the forecast loss differentials between the two competing methods equally, where the loss differential for observation t is defined by $d_t \equiv g(e_{i,t|t-h}) - g(e_{j,t|t-h})$, with $g(\cdot)$ some arbitrary loss function, and $e_{i,t|t-h}$ and $e_{j,t|t-h}$ the h -step-ahead forecast errors for Methods i and j . That is, the DM test examines whether the following equally-weighted sample mean loss differential, when standardized, is different from zero at some given significance level

$$\bar{d} = \frac{1}{P} \sum_{t=R+h}^{R+P+h-1} d_t, \quad (5)$$

where forecasts have been produced for observations $t = R+h, \dots, R+P+h-1$, such that there are P out-of-sample point forecasts and R observations have been used for estimation of the model.

Under standard conditions, Diebold and Mariano (1995) established the asymptotic normality of the DM statistic. Two important concerns with the use of DM-type statistics, however, have appeared in the literature and we address those here. First, West (1996, 2001) and West and McCracken (1998) analyzed modification of forecast comparison tests

in light of the use of estimated model parameters in the computation of such tests. van Dijk and Franses (2003) pointed out, though, that for DM-type tests under quadratic loss, such parameter estimation uncertainty is asymptotically irrelevant. They thus argued that when examining the statistical significance of MSPE reductions (which is what we are interested in), corrections of the type suggested by West (1996, 2001) and West and McCracken (1998) are not necessary.

Second, under the assumption that the estimation sample size R and the number of out-of-sample forecasts P tend to infinity, McCracken (2000) and Clark and McCracken (2001) showed that, if the underlying forecasting models are nested, the asymptotic distribution of the DM statistic is not standard normal. As noted by van Dijk and Franses (2003), these conditions on the parameters R and P effectively mean that expanding windows of data are used for estimation. In contrast, for the case in which R remains finite, as in our use of fixed-length rolling estimation windows, Giacomini and White (2003) proved that the asymptotic distribution of the DM statistic is standard normal when comparing forecasts generated by nested models.

Simulation evidence has shown that the size of the DM statistic is biased upwards in small samples. As such, Harvey *et al.* (1997) introduced a modification of the DM statistic (M-DM) to correct for this. Following Harvey *et al.* (1997), we use the Student's t distribution with $P - 1$ degrees of freedom to obtain critical values for the M-DM tests we run.

Our M-DM results appear in Table 3. Recalling that Method 10 is ranked first for the U.S. using the MSPE, we see that the MSPE of Method 10 is significantly lower, at the 10% significance level, than the MSPE of Method 11 for 50% of the forecast horizons; at no forecast step is the MSPE of Method 1 lower than that of Method 10. It is interesting to note that Method 11, which is fourth-ranked according to the MSPE and which pools across the linear and nonlinear forecasts, generates a statistically significant reduction in MSPE relative to the linear VECM case at five out of the eight forecast horizons. Investigation of these results at the individual forecast steps $h = 1, \dots, 8$ reveals that these significant MSPE reductions occur for $h = 3, \dots, 6$ for Method 10 and at $h = 3, \dots, 7$ for Method 11, that is, at the more so intermediate-term forecast horizons.

In the U.K. case, only one forecasting method, Method 10, produces a significant MSPE reduction relative to the linear VECM; and this occurs at only one forecast step. Thus, for the U.K. the linear forecasts are generally not dominated via the M-DM test at the 10% significance level. While the MSPE-top-ranked Method 1 generates a statistically significant lower MSPE relative the h -step-ahead projections of Methods 7, 8, and 9 at, respectively, six, seven, and seven forecast horizons, it does so relative to the pooling-based Methods 10 and 11 at, respectively, only two forecast steps and 1 forecast step.

For Canada, the MSPE of Methods 10 and 11 is significantly lower than that of the linear VECM-based Method 1 at, respectively, three and two forecast steps; these occur at the intermediate forecast horizons $h = 3, 5$ and 6 for the MSPE-second-ranked Method 10, and $h = 3$ and $h = 5$ for the MSPE-top-ranked Method 11. So, there is moderate evidence of statistically significant forecast improvement over the linear VECM forecasts using the two pooling procedures. On the other hand, Method 1 generates a significantly lower MSPE relative to Methods 10 and 11 at two forecast steps, the longer-term $h = 7$ and 8 .

An interesting situation is revealed by the forecasting results for Japan. While Method 1, the MSPE-top-ranked forecasting method generates a statistically significant MSPE reduction relative to the second-ranked Method 11 at only one forecast horizon, the MSPE of Method 11 is significantly lower than that of linear VECM forecasts at three forecast steps, the longer-term $h = 6, 7$ and 8 .

As a complement to the results found in Table 3, we also ran a set of Harvey *et al.* (1997) modified Diebold and Mariano (1995) tests comparing the mean absolute prediction error (MAPE) across the different forecasting methods. These results are strongly consistent with what we obtain via the M-DM test on MSPE improvements, both with respect to the frequency of significant reductions in MAPE and with respect to the particular forecast horizons at which such reductions occur.

3.3.2 Left-Tail and Right-Tail Weighting

van Dijk and Franses (2003) argued that the uniform weighting scheme employed by the M-DM test may be unsatisfactory for frequently encountered situations in which some

observations are more important than others. For example, in an unemployment forecasting exercise of the type we analyze, large positive observations for the change in the unemployment rate generally signal a business cycle downturn.

Accordingly, van Dijk and Franses modified the Diebold-Mariano statistic by weighting more heavily the loss differentials for observations that are deemed to be of greater substantive interest. In their approach, the weighted average loss differential is given by

$$\bar{d}_w = \frac{1}{P} \sum_{t=R+h}^{R+P+h-1} w(\omega_t) d_t, \quad (6)$$

where ω_t is the information set available at time t . Letting y_t be the variable to be forecast, two particular cases van Dijk and Franses studied are

$$w_{\text{LT}}(\omega_t) = 1 - \Phi(y_t), \quad (7)$$

where $\Phi(\cdot)$ is the cumulative distribution function of y_t , to focus on the left tail of the distribution of y_t , and

$$w_{\text{RT}}(\omega_t) = \Phi(y_t), \quad (8)$$

to focus on the right tail of the distribution of y_t .

A necessary condition for the associated test statistic to have an asymptotic standard normal distribution under the null hypothesis of equal forecast accuracy is that the weight function $w(\omega_t)$ be a twice continuously differentiable mapping to the $[0,1]$ interval. The weighted Diebold-Mariano statistic is computed as,

$$\text{W-DM} = \frac{\bar{d}_w}{\sqrt{\hat{V}(\bar{d}_w)}}, \quad (9)$$

where $\hat{V}(\bar{d}_w)$ is a consistent estimate of the variance of \bar{d}_w .

Following Harvey *et al.* (1997), van Dijk and Franses adjusted the W-DM statistic by way of a small-sample correction. The resulting modified W-DM statistic is given by

$$\text{MW-DM} = \sqrt{\frac{P+1-2h+h(h-1)/P}{P}} \text{W-DM}. \quad (10)$$

Once again following Harvey *et al.* (1997), van Dijk and Franses proposed using the Student's t -distribution with $P-1$ degrees of freedom to obtain critical values for the MW-DM test.

To examine the statistical significance of MSPE reductions with greater weight placed on forecast losses associated with, respectively, unemployment rate decreases and increases, we apply the left-tailed and right-tailed MW-DM tests. Comparison of the left-tailed and right-tailed results in Tables 4 and 5, both against one another and with those in Table 3, provides some interesting insight.

Consider the case of the U.S. First, using the MW-DM test with greater weight given to unemployment rate decreases, Method 10 generates significant MSPE reductions relative to the linear VECM forecasts at 75% of the forecast horizons, representing a 50% increase in comparison to uniform weighting and right-tail weighting of the forecast loss differentials. With left-tail weighting, these significant MSPE reductions occur at $h = 1, 2, 3, 4, 5,$ and $6,$ and with right-tail weighting, they occur at $h = 3, 4, 5,$ and $6.$ So, such left-tail weighting of the forecast loss differentials shows that Method 10, in addition to dominating Method 1 at the same forecast steps as with uniform and right-tail weighting, is superior to Method 1 at the short-term forecast horizons $h = 1$ and $2.$ Second, for 50% of the forecast steps with both left-tail and right-tail weighting, Method 11 generates significant MSPE reductions relative to the linear VECM forecasts. With both weighting schemes, these occur at $h = 3, 4, 5,$ and $6.$ Thus, for Method 11 versus Method 1 comparisons, use of left-tail and right-tail weighting of the forecast loss differentials produces a significant MSPE reduction at one less forecast step relative to use of uniform weighting.

For the U.K., the left-tailed MW-DM test results are, on the whole, quite similar to those obtained via uniform weighting of the forecast loss differentials. Method 11 generates a significant MSPE reduction relative to Method 1 at only one forecast step, and Method 1 significantly dominates Method 10 at a single forecast step. On the other hand, the right-tailed MW-DM test results differ considerably from the W-DM results. In particular, when the forecast loss differentials associated with unemployment rate increases are weighted more heavily, Method 11 generates significant MSPE reductions relative to the linear VECM forecasts at 50% of the forecast horizons; these occur at $h = 4, 5, 6,$ and $7.$ Also,

Method 10 dominates Method 1 at two forecast steps with right-tail weighting.

For Canada, left-tail weighting produces more significant MSPE reductions in comparing Method 11's forecasts with those of Method 1. With equal weighting, Method 11's MSPE is significantly smaller than the linear VECM's MSPE at the 10% level for two forecast steps. But with the left-tailed MW-DM test, Method 11's MSPE reductions are significant at 50% of the forecast horizons; $h = 1, 2, 3$, and 4. With unemployment rate increases weighted more heavily, Method 11's MSPE reductions relative to Method 1 are significant at the same two forecast horizons as with uniform weighting, i.e., $h = 3$ and 5.

Finally, for Japan use of the left-tailed and right-tailed MW-DM test does not generate a greater frequency of significant MSPE reductions in Method 11 versus Method 1 and Method 10 versus Method 1 comparisons. In fact, use of the MW-DM test leads to fewer such significant MSPE reductions obtained via Method 11 against the linear VECM forecasts. With uniform weighting of the forecast loss differentials, Method 11 dominates Method 1 at three forecast steps, and with left-tailed and right-tailed weighting, this occurs at, respectively, one and two forecast steps.

4 Conclusions

In this paper we set out to explore how a set of multivariate STAR models performs, both against a linear benchmark and relative to one another, in simulated real-time out-of-sample forecasting of the four non-Euro G-7 quarterly aggregate unemployment rate series. Consideration of this issue appears warranted in light of work in the empirical literature on business cycle asymmetry, in which a good deal of evidence that the data generating process for many unemployment rate series may indeed be nonlinear has been reported.

Our out-of-sample results show that, according to both forecast evaluation criteria used, the top-ranked forecasting method for the U.S. and Canada is a pooled-median forecasting approach. For the U.S., the dominant forecasting method uses the median across the set of nonlinear point forecasts; for Canada, forecasting with the median across the set of linear and nonlinear point forecasts performs best. These multivariate pooling

results are consistent with those reported by Stock and Watson (1999) and Marcellino (2004) in their analysis of univariate nonlinear models. For the U.K. and Japan, the linear VECM forecasts are top-ranked using the MSPE. Also, for three of the four countries, the relative performance of the linear forecasts worsens according to the more robust evaluation criterion; they are ranked last for the U.S. using the MedSPE criterion.

When we test the significance of the MSPE reductions we obtain using equal weighting of the forecast loss differentials, the results show that for the U.S., median-pooling across the linear and nonlinear models produces a significantly lower MSPE than that generated by the linear VECM for more than half of the forecast horizons. This statistically significant forecast improvement over the linear forecasts occurs at the more so intermediate-term forecast steps, that is, for the three-quarters-ahead to the seven-quarters-ahead forecast steps. It is interesting to see that restricting median-pooling to the nonlinear models generates fewer significant MSPE reductions over the linear VECM case, even though doing so leads to a higher ranked MSPE-ranked forecasting method. For the other countries, median-pooling leads to fewer, in comparison to the U.S. results, statistically significant decreases in MSPE relative to the linear VECM forecasts under standard uniform weighting.

When we use a recently developed test of forecast accuracy which places more weight on the forecast loss differentials associated with extreme values of the unconditional distribution of the unemployment rate first differences, some interesting behavior in the frequency of significant MSPE reductions is revealed. For the U.S., median-pooling across the nonlinear forecasts generates significant MSPE reductions over the linear VECM case 50% more often when unemployment rate decreases are emphasized as opposed to when unemployment rate increases are given more weight and to when uniform weighting is used; when unemployment rate decreases are weighted more heavily, statistically significant MSPE decreases occur at six out of the eight forecast horizons considered. Thus, these STAR forecasts perform better during business cycle expansions for the U.S. For the U.K., median-pooling across all forecast methods produces statistically significant lower MSPEs relative to the linear VECM forecasts for half of the forecast horizons when unemployment rate increases are given more weight; this occurs at short-term and intermediate-term fore-

cast steps. But when unemployment rate decreases are weighted more for the U.K., these MSPE reductions are significant at only one forecast horizon. Together these results suggest that pooling across the linear VECM and STAR forecasts works better during U.K. business cycle contractions. For Canada, application of such weighting schemes to the forecast loss differentials implies that global median-pooling across linear and nonlinear forecasts leads to more success against the linear VECM forecasts during expansionary phases of the business cycle. In contrast to the three other countries, uniform weighting of the forecast loss differentials for Japan leads to more significant MSPE reductions with median-pooling relative to what occurs when unemployment rate decreases and increases are given more weight.

We believe the main message from our forecasting exercise is as follows. While individual nonlinear forecasting methods may rarely dominate a linear approach, forecast improvement seems attainable by combining across the set of linear and nonlinear forecasts. Noting that in this paper we restrict ourselves to STAR-type multivariate models, we speculate that pooling linear forecasts with a larger set of nonlinear alternatives would prove to be useful. We intend to pursue this question in further research.

Among the set of STAR forecasting methods used, we find that no individual approach tends to outperform the others. In some cases, the top-ranked nonlinear forecasting method employs multi-step-ahead forecasts obtained by iterating the estimated one-step-ahead model. In others, h -step-ahead projections dominate. As a result, at least for the data sets examined in this paper, it appears that use of both approaches is warranted. We note that these results stand in contrast with those in Marcellino, Stock, and Watson (2004), who, in their linear study with U.S. macroeconomic time series, found that iterated forecasts generally dominated h -step-ahead projections.

In this paper we compare the point forecasts of the models used. Thus, it would be interesting to investigate the robustness of our results with respect to construction and evaluation of both interval and density forecasts. Clements and Hendry (1999, p. 285), for example, suggest that use of interval and density forecasts may indeed show improved forecasting performance for nonlinear models. We note, however, that Clements, Franses, Smith, and van Dijk (2003) report simulation results which suggest that the Diebold and

Mariano (1995) test is in fact more powerful than interval and density forecast-based tests in discriminating between linear and nonlinear models.

Table 1: Forecasting Method Definitions

Method	Definition
1	Unrestricted VECM.
2	STVECM, with transition variable selected by S_1 test run on first-differenced log-unemployment equation of subset linear VECM.
3	STVECM, with transition variable selected by S_2 test run on first-differenced log-unemployment equation of subset linear VECM.
4	STVECM, with transition variable selected by S_3 test run on first-differenced log-unemployment equation of subset linear VECM.
5	STVECM, with transition variable selected by system-wide S_1 test.
6	STVECM, with transition variable selected by system-wide S_3 test.
7	h -step-ahead projection of STVECM's first differenced log-unemployment rate equation, with transition variable selected by S_1 test run on corresponding subset equation of VECM.
8	h -step-ahead projection of STVECM's first-differenced log-unemployment rate equation, with transition variable selected by S_2 test run on corresponding subset equation of VECM.
9	h -step-ahead projection of STVECM's first-differenced log-unemployment rate equation, with transition variable selected by S_3 test run on first-differenced log-unemployment subset equation of VECM.
10	Pooled median forecast from nonlinear methods, i.e., Methods 2 through 9.
11	Pooled median forecast from Methods 1 through 9.

Table 2: Average Out-of-Sample Forecasting Ranks

Method i	U.S.	U.K.	Canada	Japan
		<u>MSPE</u>		
1	6.9	2.1	3.4	1.4
2	6.9	5.4	6.5	5.9
3	6.1	9.5	7.5	8.1
4	2.9	5.4	8.9	6.0
5	4.6	6.9	8.3	7.5
6	5.3	7.5	6.6	6.9
7	9.0	8.1	7.5	8.1
8	7.8	7.8	6.5	8.4
9	9.5	8.1	6.6	7.5
10	2.4	2.4	2.5	4.0
11	4.8	2.9	1.8	2.3
		<u>MedSPE</u>		
1	9.0	4.5	5.0	3.8
2	7.0	4.1	9.1	5.9
3	5.6	5.6	4.4	7.0
4	5.8	4.4	5.3	7.9
5	5.5	6.9	5.6	7.4
6	6.0	4.8	6.8	5.1
7	4.6	9.5	7.1	5.9
8	7.1	8.5	6.6	7.4
9	5.5	9.8	7.0	5.5
10	3.0	3.5	4.9	4.9
11	6.9	4.5	4.3	5.4

The two panels show the average out-of-sample forecasting rank of Method i across the 33 estimation windows and forecasting horizons $h = 1, \dots, 8$, using the Mean Squared Prediction Error (MSPE) and Median Squared Prediction Error (MedSPE) criteria. See Table 1 for the forecasting method definitions.

Table 3: Pair-wise Out-of-Sample Forecast Comparison Using M-DM

Method i	Method j										
	1	2	3	4	5	6	7	8	9	10	11
						<u>U.S.</u>					
1		0.0	0.0	0.0	0.0	0.0	0.0	25.0	25.0	0.0	0.0
2	12.5		0.0	0.0	0.0	0.0	25.0	50.0	12.5	0.0	0.0
3	25.0	12.5		0.0	25.0	0.0	12.5	50.0	25.0	0.0	12.5
4	50.0	37.5	25.0		37.5	25.0	25.0	50.0	37.5	12.5	50.0
5	25.0	12.5	12.5	0.0		12.5	12.5	50.0	25.0	0.0	12.5
6	25.0	25.0	0.0	12.5	25.0		25.0	50.0	25.0	0.0	25.0
7	0.0	12.5	12.5	0.0	0.0	0.0		12.5	25.0	0.0	0.0
8	12.5	25.0	12.5	0.0	0.0	12.5	0.0		25.0	0.0	0.0
9	0.0	12.5	0.0	0.0	0.0	0.0	12.5	25.0		0.0	0.0
10	50.0	50.0	12.5	12.5	62.5	50.0	37.5	50.0	50.0		50.0
11	62.5	25.0	37.5	12.5	25.0	37.5	12.5	50.0	25.0	12.5	
						<u>U.K.</u>					
1		12.5	25.0	37.5	12.5	0.0	75.0	87.5	87.5	25.0	12.5
2	0.0		12.5	0.0	0.0	0.0	50.0	50.0	50.0	0.0	0.0
3	0.0	0.0		0.0	0.0	0.0	37.5	12.5	25.0	0.0	0.0
4	0.0	0.0	0.0		0.0	0.0	50.0	37.5	50.0	0.0	0.0
5	0.0	0.0	25.0	0.0		0.0	25.0	37.5	37.5	0.0	12.5
6	0.0	0.0	12.5	0.0	12.5		25.0	25.0	37.5	0.0	0.0
7	0.0	0.0	0.0	0.0	0.0	0.0		0.0	0.0	0.0	0.0
8	0.0	0.0	0.0	0.0	0.0	0.0	0.0		0.0	0.0	0.0
9	0.0	0.0	0.0	0.0	0.0	0.0	12.5	0.0		0.0	0.0
10	12.5	25.0	37.5	37.5	12.5	0.0	87.5	100.0	87.5		25.0
11	0.0	25.0	0.0	62.5	12.5	0.0	87.5	100.0	87.5	50.0	

continued on next page

continued from previous page

Method i	Method j										
	1	2	3	4	5	6	7	8	9	10	11
	<u>Canada</u>										
1		62.5	37.5	12.5	0.0	12.5	50.0	37.5	37.5	25.0	25.0
2	0.0		0.0	0.0	0.0	0.0	12.5	12.5	12.5	0.0	0.0
3	0.0	0.0		12.5	0.0	0.0	0.0	0.0	0.0	0.0	0.0
4	0.0	0.0	0.0		0.0	0.0	0.0	0.0	0.0	0.0	0.0
5	0.0	0.0	12.5	12.5		0.0	0.0	0.0	0.0	0.0	0.0
6	0.0	12.5	0.0	12.5	12.5		12.5	12.5	0.0	0.0	0.0
7	12.5	12.5	12.5	12.5	12.5	12.5		12.5	12.5	12.5	0.0
8	12.5	12.5	12.5	12.5	12.5	12.5	12.5		25.0	12.5	0.0
9	12.5	12.5	12.5	12.5	12.5	12.5	12.5	0.0		12.5	0.0
10	37.5	62.5	37.5	12.5	25.0	37.5	50.0	50.0	37.5		0.0
11	25.0	62.5	62.5	12.5	37.5	25.0	50.0	62.5	50.0	12.5	
	<u>Japan</u>										
1		0.0	25.0	37.5	50.0	12.5	25.0	37.5	25.0	25.0	12.5
2	0.0		25.0	37.5	50.0	12.5	12.5	0.0	12.5	12.5	0.0
3	0.0	0.0		25.0	0.0	0.0	12.5	0.0	0.0	0.0	0.0
4	0.0	12.5	25.0		25.0	0.0	0.0	0.0	0.0	0.0	0.0
5	0.0	12.5	12.5	0.0		0.0	12.5	0.0	0.0	0.0	0.0
6	0.0	0.0	0.0	12.5	37.5		12.5	0.0	12.5	0.0	0.0
7	0.0	0.0	0.0	0.0	12.5	0.0		12.5	12.5	0.0	0.0
8	0.0	0.0	0.0	0.0	12.5	0.0	0.0		0.0	0.0	0.0
9	0.0	0.0	0.0	0.0	12.5	0.0	37.5	0.0		0.0	0.0
10	0.0	37.5	50.0	25.0	50.0	25.0	25.0	0.0	12.5		0.0
11	37.5	25.0	25.0	37.5	37.5	50.0	12.5	25.0	25.0	25.0	

The table presents pair-wise out-of-sample forecast comparisons for the 11 forecasting methods and 33 estimation windows, across forecasting horizons $h = 1, \dots, 8$, using the modified Diebold-Mariano MSPE statistic of Harvey *et al.* (1997) (M-DM). The entries in the table show the percentage of forecast horizons for which the M-DM test rejects the null hypothesis that Method i 's forecast performance as measured by MSPE is not superior to that of Method j at the 10% significance level. See Table 1 for the forecasting method definitions.

Table 4: Pair-wise Out-of-Sample Forecast Comparison Using Left-Tailed MW-DM

Method i	Method j										
	1	2	3	4	5	6	7	8	9	10	11
						<u>U.S.</u>					
1		25.0	12.5	0.0	0.0	0.0	0.0	12.5	12.5	0.0	0.0
2	50.0		12.5	0.0	0.0	0.0	0.0	50.0	12.5	0.0	25.0
3	37.5	12.5		0.0	25.0	25.0	0.0	37.5	12.5	0.0	37.5
4	75.0	25.0	12.5		25.0	25.0	0.0	25.0	12.5	0.0	37.5
5	50.0	25.0	0.0	0.0		25.0	0.0	25.0	12.5	0.0	25.0
6	37.5	12.5	0.0	0.0	25.0		0.0	37.5	12.5	0.0	25.0
7	12.5	12.5	0.0	12.5	0.0	12.5		12.5	12.5	0.0	0.0
8	12.5	25.0	12.5	12.5	12.5	12.5	0.0		25.0	0.0	0.0
9	0.0	0.0	0.0	0.0	0.0	0.0	12.5	12.5		0.0	0.0
10	75.0	62.5	25.0	25.0	62.5	50.0	25.0	50.0	12.5		50.0
11	50.0	37.5	25.0	12.5	50.0	37.5	0.0	12.5	12.5	12.5	
						<u>U.K.</u>					
1		37.5	12.5	50.0	37.5	12.5	75.0	87.5	75.0	12.5	0.0
2	0.0		0.0	0.0	0.0	0.0	62.5	50.0	62.5	0.0	0.0
3	0.0	0.0		0.0	0.0	0.0	25.0	25.0	25.0	0.0	0.0
4	0.0	12.5	0.0		0.0	0.0	62.5	62.5	62.5	0.0	0.0
5	0.0	0.0	37.5	12.5		0.0	62.5	50.0	62.5	0.0	0.0
6	0.0	0.0	0.0	0.0	37.5		50.0	37.5	50.0	0.0	0.0
7	0.0	12.5	0.0	0.0	0.0	0.0		12.5	0.0	0.0	0.0
8	0.0	0.0	0.0	0.0	0.0	0.0	12.5		12.5	0.0	0.0
9	0.0	12.5	0.0	0.0	0.0	0.0	0.0	12.5		0.0	0.0
10	0.0	12.5	37.5	12.5	37.5	0.0	75.0	100.0	87.5		12.5
11	12.5	37.5	12.5	37.5	62.5	25.0	75.0	100.0	87.5	50.0	

continued on next page

continued from previous page

Method i	Method j										
	1	2	3	4	5	6	7	8	9	10	11
<u>Canada</u>											
1		37.5	37.5	0.0	0.0	0.0	37.5	37.5	50.0	25.0	25.0
2	0.0		0.0	0.0	0.0	0.0	12.5	12.5	12.5	0.0	0.0
3	0.0	0.0		12.5	0.0	0.0	0.0	0.0	0.0	0.0	0.0
4	0.0	0.0	0.0		12.5	0.0	0.0	0.0	0.0	0.0	0.0
5	0.0	0.0	0.0	0.0		12.5	0.0	0.0	0.0	0.0	0.0
6	0.0	12.5	0.0	0.0	12.5		0.0	0.0	0.0	0.0	0.0
7	12.5	12.5	25.0	12.5	12.5	12.5		12.5	0.0	0.0	0.0
8	12.5	12.5	12.5	12.5	12.5	12.5	12.5		12.5	0.0	0.0
9	12.5	12.5	12.5	12.5	12.5	12.5	12.5	25.0		12.5	0.0
10	12.5	75.0	37.5	12.5	25.0	37.5	50.0	50.0	62.5		0.0
11	50.0	75.0	62.5	12.5	25.0	25.0	62.5	62.5	50.0	25.0	
<u>Japan</u>											
1		37.5	50.0	87.5	87.5	37.5	37.5	75.0	37.5	37.5	0.0
2	0.0		25.0	25.0	50.0	12.5	12.5	12.5	12.5	12.5	0.0
3	0.0	0.0		12.5	0.0	0.0	12.5	12.5	12.5	0.0	0.0
4	0.0	0.0	0.0		25.0	0.0	12.5	12.5	12.5	0.0	0.0
5	0.0	0.0	0.0	12.5		0.0	12.5	0.0	12.5	0.0	0.0
6	0.0	0.0	0.0	12.5	25.0		12.5	25.0	12.5	0.0	0.0
7	0.0	0.0	0.0	0.0	0.0	0.0		12.5	12.5	0.0	0.0
8	0.0	0.0	0.0	0.0	0.0	0.0	0.0		12.5	0.0	0.0
9	0.0	0.0	0.0	0.0	0.0	0.0	25.0	12.5		0.0	0.0
10	0.0	37.5	12.5	25.0	62.5	25.0	12.5	37.5	12.5		0.0
11	12.5	50.0	62.5	87.5	100.0	50.0	50.0	75.0	50.0	50.0	

The table presents pair-wise out-of-sample forecast comparisons for the 11 forecasting methods and 33 estimation windows, across forecasting horizons $h = 1, \dots, 8$, using the left-tailed modified weighted Diebold-Mariano MSPE statistic of van Dijk and Franses (2003) (MW-DM). The entries in the table show the percentage of forecast horizons for which the left-tailed MW-DM test rejects the null hypothesis that Method i 's forecast performance as measured by MSPE is not superior to that of Method j at the 10% significance level. See Table 1 for the forecasting method definitions.

Table 5: Pair-wise Out-of-Sample Forecast Comparison Using Right-Tailed MW-DM

Method i	Method j										
	1	2	3	4	5	6	7	8	9	10	11
						<u>U.S.</u>					
1		0.0	12.5	0.0	0.0	0.0	12.5	25.0	37.5	0.0	0.0
2	0.0		0.0	0.0	0.0	0.0	25.0	62.5	25.0	0.0	0.0
3	37.5	12.5		0.0	25.0	12.5	12.5	50.0	50.0	0.0	12.5
4	50.0	37.5	37.5		50.0	37.5	62.5	62.5	50.0	12.5	25.0
5	25.0	25.0	12.5	0.0		25.0	37.5	62.5	50.0	0.0	0.0
6	12.5	37.5	0.0	0.0	37.5		37.5	62.5	62.5	0.0	12.5
7	0.0	0.0	0.0	0.0	0.0	0.0		25.0	12.5	0.0	0.0
8	0.0	0.0	12.5	0.0	0.0	0.0	0.0		12.5	0.0	0.0
9	0.0	12.5	12.5	0.0	0.0	0.0	12.5	25.0		0.0	0.0
10	50.0	37.5	12.5	0.0	50.0	37.5	50.0	62.5	62.5		12.5
11	50.0	12.5	25.0	0.0	12.5	12.5	50.0	62.5	50.0	0.0	
						<u>U.K.</u>					
1		0.0	25.0	37.5	0.0	0.0	75.0	62.5	75.0	0.0	12.5
2	0.0		12.5	0.0	0.0	0.0	37.5	50.0	37.5	0.0	0.0
3	0.0	0.0		0.0	0.0	0.0	37.5	12.5	37.5	0.0	0.0
4	0.0	0.0	12.5		0.0	0.0	37.5	37.5	37.5	0.0	0.0
5	0.0	0.0	25.0	0.0		0.0	25.0	25.0	25.0	0.0	12.5
6	0.0	0.0	12.5	0.0	25.0		25.0	12.5	25.0	0.0	12.5
7	0.0	0.0	0.0	0.0	0.0	0.0		12.5	12.5	0.0	0.0
8	0.0	0.0	0.0	0.0	0.0	0.0	12.5		12.5	0.0	0.0
9	0.0	0.0	0.0	0.0	0.0	0.0	12.5	0.0		0.0	0.0
10	25.0	12.5	25.0	37.5	0.0	12.5	87.5	75.0	75.0		25.0
11	50.0	12.5	0.0	37.5	0.0	0.0	75.0	62.5	75.0	0.0	

continued on next page

continued from previous page

Method i	Method j										
	1	2	3	4	5	6	7	8	9	10	11
	<u>Canada</u>										
1		37.5	37.5	25.0	12.5	12.5	37.5	37.5	25.0	0.0	25.0
2	0.0		0.0	0.0	0.0	0.0	25.0	12.5	0.0	0.0	0.0
3	0.0	12.5		0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
4	0.0	12.5	0.0		12.5	0.0	0.0	0.0	0.0	0.0	0.0
5	0.0	0.0	0.0	12.5		0.0	0.0	0.0	0.0	0.0	0.0
6	0.0	12.5	0.0	25.0	25.0		25.0	12.5	0.0	0.0	0.0
7	0.0	25.0	0.0	0.0	12.5	12.5		12.5	12.5	0.0	0.0
8	0.0	25.0	0.0	12.5	12.5	12.5	0.0		12.5	0.0	0.0
9	0.0	12.5	0.0	0.0	12.5	12.5	12.5	0.0		0.0	0.0
10	12.5	62.5	12.5	12.5	25.0	25.0	37.5	50.0	37.5		12.5
11	25.0	37.5	25.0	37.5	50.0	37.5	37.5	37.5	50.0	12.5	
	<u>Japan</u>										
1		0.0	25.0	12.5	12.5	12.5	25.0	12.5	12.5	12.5	37.5
2	0.0		25.0	37.5	37.5	0.0	0.0	0.0	0.0	0.0	0.0
3	0.0	0.0		25.0	0.0	12.5	0.0	0.0	0.0	0.0	0.0
4	0.0	12.5	25.0		25.0	0.0	0.0	0.0	0.0	0.0	12.5
5	0.0	12.5	12.5	0.0		0.0	0.0	0.0	0.0	0.0	0.0
6	0.0	0.0	25.0	0.0	37.5		12.5	0.0	12.5	0.0	0.0
7	0.0	0.0	0.0	12.5	12.5	0.0		25.0	0.0	0.0	0.0
8	0.0	0.0	0.0	0.0	12.5	12.5	0.0		0.0	0.0	0.0
9	0.0	0.0	0.0	12.5	12.5	0.0	12.5	12.5		0.0	0.0
10	0.0	37.5	37.5	12.5	50.0	12.5	12.5	0.0	0.0		12.5
11	25.0	25.0	12.5	25.0	25.0	25.0	12.5	12.5	12.5	12.5	

The table presents pair-wise out-of-sample forecast comparisons for the 11 forecasting methods and 33 estimation windows, across forecasting horizons $h = 1, \dots, 8$, using the right-tailed modified weighted Diebold-Mariano MSPE statistic of van Dijk and Franses (2003) (MW-DM). The entries in the table show the percentage of forecast horizons for which the right-tailed MW-DM test rejects the null hypothesis that Method i 's forecast performance as measured by MSPE is not superior to that of Method j at the 10% significance level. See Table 1 for the forecasting method definitions.

Endnotes

¹See Table 11 of Marcellino (2002).

²Other papers which study the performance of nonlinear time series models in forecasting unemployment rate fluctuations include Peel and Speight (2000), Terui and van Dijk (2002), and Proietti (2003).

³In both Rothman *et al.* (2001) and here, revised as opposed to real-time or preliminary data are used; see, for example, Amato and Swanson (2001) for discussion of the distinction between these. The main reason we use revised data is that real-time data sets for all of the systems we estimate are not available. While this certainly warrants the standard caveat about our results, what we do is also consistent with, for example, Stock and Watson (1999) and Stock and Watson (2001).

⁴In this paper, which appears in a special issue of *Empirical Economics* devoted to recent developments in modelling business cycle and financial data with regime-switching models, the one-step-ahead forecasts were computed for a given parametric structure via an updating of the regime-dependent probabilities.

⁵We decided to work with log-unemployment rate data following the theoretical framework presented in Nickell (1998) and in order to reduce the heteroskedasticity of the residuals in our estimated models.

⁶The money supply and interest rate series used are M2 and the 90-day Treasury bill rate for the U.S., M4 and the 90-day Treasury bill rate for the U.K., M2 and the 90-day commercial paper rate for Canada, and M2 and the lending rate for collateral and overnight loans in the Tokyo call money market for Japan. We use M4 data for the UK since the M2 time series is incomplete and inconsistent. We use the overnight Tokyo call rate for Japan since no sufficient long 90-day rate is available; see, for example, Table 1 of Stock and Watson (2001)

⁷This procedure leads to the model that would be selected by applying the AIC to each equation individually. But it is not guaranteed that this model also minimizes the AIC for the system as a whole. The simulation evidence in Brüggemann and Lütkepohl (2000), however, shows that the difference between the models selected by this single equation approach and a comparable system approach is generally small.

⁸We found, by way of generating an alternative set of forecasts with u_t removed from \mathbf{z}_t , that our results are robust with respect to the assumption that the unemployment rate is $I(1)$ instead of $I(0)$.

⁹More details on these system-wide versions of the Luukkonen *et al.* (1988) tests can be found in Rothman *et al.* (2001).

¹⁰The reason why we do not use analogues of the federal funds rate for the U.K. and Canada is that, if we were to do so, this would shorten considerably the available time series; see, once again, Table 1 of Stock and Watson (2001). Also, as noted earlier, our short-term interest rate for Japan is an overnight rate.

¹¹Details on all top-ranked transition variables used in this paper, along with p -values of

the associated linearity tests, are available upon request. Since our primary interest in this paper is on out-of-sample forecasting, we do not focus on the linearity testing results here. That said, we note that the battery of linearity tests run reveal strong evidence in favor of STAR-type nonlinearity and that the rankings of the candidate transition variables vary a good deal across the particular tests employed and the individual unemployment rate series examined.

¹²These are available upon request.

References

- Altissimo, F. and G.L. Violante (2001), “The Non-Linear Dynamics of Output and Unemployment in the U.S.,” *Journal of Applied Econometrics* 16, 461–486.
- Amato, J.D. and N.R. Swanson (2001), “The Real-Time Predictive Content of Money for Output,” *Journal of Monetary Economics* 48, 3–24.
- Balke, N.S. and T.B. Fomby (1997), “Threshold Cointegration,” *International Economic Review* 38, 627–646.
- Beaudry, P. and G. Koop (1993), “Do Recessions Permanently Affect Output?” *Journal of Monetary Economics* 31, 149–163.
- Brüggemann, R. and H. Lütkepohl (2000), “Lag Selection in Subset Var Models with an Application to a U.S. Monetary System,” unpublished manuscript, Humboldt University.
- Caballero, R. and M. Hammour (1994), “The Cleansing Effects of Recession,” *American Economic Review* 84, 1350–1368.
- Caner, M. and B.E. Hansen (2001), “Threshold Autoregression with a Unit Root,” *Econometrica* 69, 1555–1596.
- Caplin, A.S. and J. Leahy (1991), “State-Contingent Pricing and the Dynamics of Money and Output,” *Quarterly Journal of Economics* 106, 683–708.
- Clark, T. E. and M.W. McCracken (2001), “Tests of Equal Forecast Accuracy and Encompassing for Nested Models,” *Journal of Econometrics* 105, 85–110.
- Clements, M.P., P.H. Franses, J. Smith, and D. van Dijk (2003), “On SETAR Non-Linearity and Forecasting,” *Journal of Forecasting* 22, 359–375.
- Clements, M.P. and D.H. Hendry (1999), *Forecasting Non-stationary Economic Time Series*, Cambridge, MA: MIT Press.
- Clements, M.P. and H.-M. Krolzig (2003), “Business Cycle Asymmetries: Characterization and Testing Based on Markov-Switching Autoregressions,” *Journal of Business & Economic Statistics* 21, 196–211.
- Cooley, T.F. and G.D. Hansen (1995), “Money and the Business Cycle,” in Cooley, T.F., (ed.), *Frontiers of Business Cycle Research*, Princeton, NJ: Princeton University Press, 175–216.
- Davies, R.B. (1987), “Hypothesis Testing When a Nuisance Parameter is Present Only Under the Alternative,” *Biometrika* 74, 33–43.
- Diebold, F.X. and R.S. Mariano (1995), “Comparing Predictive Accuracy,” *Journal of Business and Economic Statistics* 13, 253–263.
- Giacomini, R. and White, H. (2003), “Tests of Conditional Predictive Ability,” UCSD Working Paper No. 2003-09, University of California, San Diego.
- Harvey, D.I., S.J. Leybourne, and P. Newbold (1997), “Testing the Equality of Prediction Mean Squared Errors,” *International Journal of Forecasting* 13, 281–291.

- Hendry, D.F. and G.E. Mizon (1993), "Evaluating Dynamic Econometric Models by Encompassing the VAR," in Phillips, P.C.B. (ed.), *Models, Methods, and Applications: Essays in Honor of A.R. Bergstrom*, Cambridge, MA: Basil Blackwell, 272–300.
- Krolzig, H.-M., M. Marcellino, and G.E. Mizon (2002), "A Markov-Switching Vector Equilibrium Correction Model of the UK Labour Market," *Empirical Economics* 27, 233–254.
- Lundbergh, S. and T. Teräsvirta (1998), "Modelling Economic High-Frequency Time Series with STAR-GARCH Models," Working Paper Series in Economics and Finance No. 291, Stockholm School of Economics.
- Luukkonen, R., P. Saikkonen and T. Teräsvirta (1988), "Testing Linearity Against Smooth Transition Autoregressive Models," *Biometrika* 75, 491–499.
- Marcellino, M. (2002), "Instability and Non-Linearity in the EMU," Centre for Economic Policy Research Discussion Paper No. 3312.
- Marcellino, M. (2004), "Forecast Pooling for Short Time Series of Macroeconomic Variables," *Oxford Bulletin of Economics & Statistics* 66, 91–112.
- Marcellino, M., J.H. Stock, and M.W. Watson (2004), "A Comparison of Direct and Iterated Multistep AR Methods for Forecasting Macroeconomic Time Series," unpublished manuscript, Università Bocconi.
- McCracken, M. W. (2000), "Robust Out-of-Sample Inference," *Journal of Econometrics* 99, 195–223.
- Montgomery, A.L., V. Zarnowitz, R.S. Tsay, and G.C. Tiao (1998), "Forecasting the U.S. Unemployment Rate," *Journal of the American Statistical Association* 93, 478–493.
- Neftci, S.N. (1984), "Are Economic Time Series Asymmetric Over the Business Cycle?" *Journal of Political Economy* 92, 307–328.
- Nickell, S. (1998), "Unemployment: Questions and Some Answers," *The Economic Journal* 108, 802–816.
- Peel, D.A. and A.E.H. Speight (2000), "Threshold Nonlinearities in Unemployment Rates: Further Evidence for the UK and G3 Economies," *Applied Economics* 32, 705–715.
- Proietti, T. (2003), "Forecasting the US Unemployment Rate," *Computational Statistics & Data Analysis* 42, 451–476.
- Rothman, P. (1998), "Forecasting Asymmetric Unemployment Rates," *Review of Economics & Statistics* 80, 164–168.
- Rothman, P., D. van Dijk, and P.H. Franses (2001), "A Multivariate STAR Analysis of the Relationship Between Money and Output," *Macroeconomic Dynamics* 5, 506–532.
- Skalin, J. and T. Teräsvirta (2002), "Modeling Asymmetries and Moving Equilibria in Unemployment Rates," *Macroeconomic Dynamics* 6, 202–241.
- Söderlind, P. and A. Vredin (1996), "Applied Cointegration Analysis in the Mirror of Macroeconomic Theory," *Journal of Applied Econometrics* 11, 363–381.

- Stock, J.H. and M.W. Watson (2001), "Forecasting Output and Inflation: The Role of Asset Prices," NBER Working Paper No. 8180.
- Stock, J.H. and M.W. Watson (1999), "A Comparison of Linear and Nonlinear Univariate Models for Forecasting Macroeconomic Time Series," in Engle, R. and H. White (eds.), *Cointegration, Causality, and Forecasting: A Festschrift in Honor of Clive W.J. Granger*, Oxford: Oxford University Press, 1–44.
- Teräsvirta, T. (1994), "Specification, Estimation, and Evaluation of Smooth Transition Autoregressive Models," *Journal of the American Statistical Association* 89, 208–218.
- Terui, N. and H.K. van Dijk (2002), "Combined Forecasts from Linear and Nonlinear Time Series Models," *International Journal of Forecasting* 18, 421–438.
- Tsay, R.S. (1998), "Testing and Modeling Multivariate Threshold Models," *Journal of the American Statistical Association* 93, 1188–1202.
- van Dijk, D. and P.H. Franses (2003), "Selecting a Nonlinear Time Series Model Using Weighted Tests of Equal Forecast Accuracy," *Oxford Bulletin of Economics & Statistics* 65, 727–744.
- West, K. D. (2001), "Tests for Forecast Encompassing When Forecasts Depend on Estimated Regression Parameters," *Journal of Business & Economic Statistics* 19, 29–33.
- West, K. D. (1996), "Asymptotic Inference About Predictive Ability," *Econometrica* 64, 1067–1084.
- West, K. D. and M.W. McCracken (1998), "Regression Based Tests of Predictive Ability," *International Economic Review* 39, 817–840.
- Wooldridge, J.M. (1991), "On the Application of Robust, Regression-Based Diagnostics to Models of Conditional Means and Conditional Variance," *Journal of Econometrics* 47, 5–46.