

Path Forecast Evaluation

Abstract

A forecast path refers to the vector of forecasts over the next 1 to h periods into the future. These forecasts are correlated across horizons so that to properly understand the uncertainty associated with the forecast path, one requires the joint predictive density of the path rather than the collection of marginal predictive densities for each horizon. This paper derives the joint predictive density for forecasts generated by VARs or from direct forecast methods (e.g. Marcellino, Stock and Watson, 2003). Given this density, we introduce the mean square forecast path metric to compare the predictive ability between competing models and appropriately modify Diebold-Mariano-West and Giacomini-White predictive ability tests. We then use Scheffé's S-method to construct simultaneous confidence regions for the forecast path and show how to construct path forecasts conditional on assumed paths for a subset of the system's variables, along with their conditional predictive density and a test on the assumed path's likelihood.

JEL Classification Codes: C32, C52, C53

Keywords: path forecast, local projections, vector autoregression, simultaneous confidence region.

Òscar Jordà
Department of Economics
University of California, Davis
One Shields Ave.
Davis, CA 95616-8578
e-mail: ojorda@ucdavis.edu

Massimiliano Marcellino
Università Bocconi, IGIER and CEPR
Via Salasco 5
20136 Milan, Italy
e-mail: massimiliano.marcellino@uni-bocconi.it

1 Introduction

Understanding the uncertainty associated with a forecast is as important as the forecast itself. When predictions are made over several periods, such uncertainty is encapsulated by the joint density of the *forecast path*. There are many questions of interest that can be answered based on the marginal distribution of the forecasts at each individual horizon. These are the questions that have received the bulk of attention in the literature and are coded into most commercial econometric packages. For example, mean-squared forecast errors (MSFE) are reported for each forecast horizon individually; two standard-error band plots that are based on the marginal distribution of each individual forecast error; and fan charts that are constructed from the percentiles of marginal predictive densities.

The basic message of this paper is that many questions of interest require knowledge of the joint density, not the collection of marginal densities alone. The joint distribution and the covariance matrix for the forecast path thus play a prominent role in our discussion, and we begin by deriving appropriate asymptotic results for data generating processes (DGP) of infinite order. Vector autoregressions are a natural multivariate method of producing forecasts and we will provide results that complement those available in, e.g., Lütkepohl (2005). Alternatively, direct forecast methods (see, e.g., Bhansali, 2002 and references therein, and more recently Marcellino, Stock and Watson, 2003) are a natural choice when there is hesitation about the true model characterizing the DGP or when nonlinearities make multiple-step ahead forecasts cumbersome to obtain. We derive asymptotic results for direct forecasts based on linear vector autoregressions for infinite order DGPs. We will call forecasts based on this method *local projections*, following the nomenclature in Jordà (2005). Some of the

derivations that we provide will look familiar to those readers acquainted with Lewis and Reinsel (1985), and Kuersteiner (2001, 2002), to cite a few.

A natural consequence of the correlation across forecast horizons is the desire to cast forecasting performance comparisons in terms of forecast paths. Therefore, we introduce the mean squared forecast path (MSFP) as the natural extension of the MSFE by using the Wald metric and the forecast path's correlation matrix. Furthermore, we show how to appropriately use this Wald metric to extend formal testing of predictive ability along the lines of the Diebold-Mariano-West (Diebold and Mariano, 1995; West, 1996) test or the more recent approach in Giacomini and White (2006). These extensions also complement deterministic measures such as the determinant of the covariance matrix of the vector of forecast errors at different horizons proposed by Clements and Hendry (1993).

A 95% confidence, multi-dimensional ellipse based on the joint distribution of the forecast path is an accurate representation of its uncertainty but is impossible to display in two-dimensional space. Another contribution of the paper is to introduce several methods to present such joint uncertainty in a useful manner to the end-user of the forecasting exercise based on Scheffé's (1953) method of simultaneous inference. In particular, we show how to construct simultaneous confidence bands, conditional confidence bands for the uncertainty associated to individual forecast horizon, and fan charts based on the quantiles of the joint predictive density.

The availability of the joint predictive density allows us to construct the distribution of forecasts conditional on future values of one or more of the endogenous variables in the system under consideration. The Wald metric provides a natural statistic to evaluate the likelihood of

observing the conditioning paths. These results can be thought of as the large-sample versions of Waggoner and Zha’s (1999) Bayesian derivations and provide asymptotic justification for bootstrap-based, finite-sample inference (Horowitz, 2001). We provide empirical examples to illustrate all of the results introduced in the paper.

Several final comments are worth making. We hope our paper will pave the way for many natural extensions that: (1) generalize our basic assumptions on the DGP as well as the mixing and heteroskedasticity assumptions of the error process; and (2) extend our large-sample results to finite-sample inference based on bootstrap or subsampling refinements. In the end, we had to allow space considerations and transparency to rein in our ambitions. Finally, ours is not a criticism of the status-quo, rather we view our contribution as an addition to existing methods: different hypotheses require different statistics and tailoring the statistics results in more precise answers.

2 Asymptotic Distribution of the Forecast Path

This section characterizes the asymptotic distribution of the forecast path, under the assumption that the data generating process (DGP) is of infinite order while the forecasts are generated by finite-order VARs or finite-order local projections. We feel the DGP is sufficiently general to represent a large class of problems of interest and that VARs and local projections are the two most commonly used modeling strategies. We begin by stating our assumptions on the DGP.

Assumption 1: Suppose the k -dimensional vector of weakly stationary variables, \mathbf{y}_t has

a Wold representation given by

$$\mathbf{y}_t = \boldsymbol{\mu} + \sum_{j=0}^{\infty} \Phi_j' \mathbf{u}_{t-j}, \quad (1)$$

where the moving-average coefficient matrices Φ_j are of dimension $k \times k$, and we assume that:

- (i) $E(\mathbf{u}_t) = 0$; and u_t are *i.i.d.*
- (ii) $E(\mathbf{u}_t \mathbf{u}_t') = \Sigma_u < \infty$.
- (iii) $\sum_{j=0}^{\infty} \|\Phi_j\| < \infty$ where $\|\Phi_j\|^2 = \text{tr}(\Phi_j' \Phi_j)$ is the equivalent of the Euclidean L_2 norm for matrices and $\Phi_0 = I_k$.
- (iv) $\det\{\Phi(z)\} \neq 0$ for $|z| \leq 1$ where $\Phi(z) = \sum_{j=0}^{\infty} \Phi_j z^j$.

Then the process in (1) can also be written as an infinite VAR process (see, e.g. Anderson, 1994),

$$\mathbf{y}_t = \mathbf{m} + \sum_{j=1}^{\infty} A_j \mathbf{y}_{t-j} + \mathbf{u}_t \quad (2)$$

such that,

- (v) $\sum_{j=1}^{\infty} \|A_j\| < \infty$.
- (vi) $A(z) = I_k - \sum_{j=1}^{\infty} A_j z^j = \Phi(z)^{-1}$.
- (vii) $\det\{A(z)\} \neq 0$ for $|z| \leq 1$.

Assumption 1 includes the class of stationary vector autoregressive moving average, VARMA(p, q) processes as a special case. Lewis and Reinsel (1985) derive conditions under

which a finite order VAR will provide consistent and asymptotically normal estimates of the p original autoregressive coefficient matrices A_j in expression (2). We will use this result momentarily and extend it for local projections when deriving the asymptotic distribution of the forecast path. The *i.i.d.* assumption could be relaxed to allow for heteroskedasticity so that the consistency and asymptotic normality results in Lewis and Reinsel (1985) are derived with appropriate laws of large numbers and central limit theorems for martingale difference sequences (*m.d.s.*) under more general mixing conditions. We refer the reader to Gonçalves and Kilian (2006) and references therein for a discussion of these issues. The most significant implication of allowing for these alternative, more flexible assumptions is that a robust covariance estimate along the lines of White (1980) is advised. For now, we prefer to trade-off some sophistication for clarity to illustrate the more important points we discuss below.

Given the DGP in expression (2) suppose we estimate a VAR(p) instead. This VAR(p) will be of the form

$$\begin{aligned} \mathbf{y}_t &= \mathbf{m} + \sum_{j=1}^p A_j \mathbf{y}_{t-j} + \mathbf{w}_t \\ \mathbf{w}_t &= \sum_{j=p+1}^{\infty} A_j \mathbf{y}_{t-j} + \mathbf{u}_t \end{aligned} \tag{3}$$

Given estimates from this VAR(p) then one could construct forecasts with standard available formulas (see, e.g. Hamilton, 1994). Alternatively, forecasts could be constructed with a sequence of local projections given by

$$\begin{aligned}
\mathbf{y}_{t+h} &= \mathbf{m}_h + \sum_{j=0}^{p-1} A_j^h \mathbf{y}_{t-j} + \mathbf{v}_{t+h} \\
\mathbf{v}_{t+h} &= \sum_{j=p}^{\infty} A_j^h \mathbf{y}_{t-j} + \mathbf{u}_{t+h} + \sum_{j=1}^{h-1} \Phi_j \mathbf{u}_{t+h-j} \quad \text{for } h = 1, \dots, H
\end{aligned} \tag{4}$$

where:

$$(i) \quad A_1^h = \Phi_h \text{ for } h \geq 1$$

$$(ii) \quad A_j^h = \Phi_{h-1} A_j + A_{j+1}^{h-1} \text{ for } h \geq 1; A_{j+1}^0 = 0; \Phi_0 = I_k; \text{ and } j \geq 1$$

Let $\Gamma(j) \equiv E(\mathbf{y}_t \mathbf{y}_{t+j}')$ with $\Gamma(-j) = \Gamma(j)'$ and define:

$$(iii) \quad X_{t,p} = (\mathbf{1}, \mathbf{y}_{t-1}', \dots, \mathbf{y}_{t-p}')'.$$

$$(iv) \quad \widehat{\Gamma}_{1-p,h} = (T-p-h)^{-1} \sum_{t=p}^T X_{t,p} \mathbf{y}_{t+h}'.$$

$$(v) \quad \widehat{\Gamma}_p = (T-p-h)^{-1} \sum_{t=p}^T X_{t,p} X_{t,p}'.$$

Then, the least-squares estimate of the VAR(p) in expression (3) is given by the formula

$$\widehat{A}(p)_{k \times kp+1} = (\widehat{\mathbf{m}}, \widehat{A}_1, \dots, \widehat{A}_p) = \widehat{\Gamma}_{1-p,0}' \widehat{\Gamma}_p^{-1},$$

whereas the coefficients of the mean-squared error linear predictor of \mathbf{y}_{t+h} based on $\mathbf{y}_t, \dots, \mathbf{y}_{t-p+1}$

is given by the least-squares formula

$$\widehat{A}(p, h)_{k \times kp+1} = (\widehat{\mathbf{m}}_h, \widehat{A}_1^h, \dots, \widehat{A}_p^h) = \widehat{\Gamma}_{1-p,h}' \widehat{\Gamma}_p^{-1}; \quad h = 1, \dots, H.$$

Assumption 2: If $\{\mathbf{y}_t\}$ satisfies conditions (i)-(vii) in assumption 1 and:

(i) $E |u_{it}u_{jt}u_{rt}u_{lt}| < \infty$ for $1 \leq i, j, r, l \leq k$.

(ii) p is chosen as a function of T such that

$$\frac{p^3}{T} \rightarrow 0 \text{ as } T, p \rightarrow \infty.$$

(iii) p is chosen as a function of T such that

$$p^{1/2} \sum_{j=p+1}^{\infty} ||A_j|| \rightarrow 0 \text{ as } T, p \rightarrow \infty.$$

Then, a summary of results shown by Lewis and Reinsel (1985), Lütkepohl and Poskitt (1991) and Jordà and Kozicki (2007) are contained in the following corollary.

Corollary 1 *Given assumptions 1 and 2, the VAR(p) and p^{th} order local projections are consistent and asymptotically normal, specifically:*

(a) $\hat{A}_j \xrightarrow{p} A_j$; $\hat{A}_j^h \xrightarrow{p} A_j^h$ and $\hat{A}_1^h \xrightarrow{p} \Phi_h$.

(b) $\sqrt{\frac{T-p-h}{p}} \text{vec} \left(\hat{A}(p) - A(p) \right) \xrightarrow{d} N(0, \Sigma_a)$ where $\Sigma_a = \Gamma_p^{-1} \otimes \Sigma_u$

(c) $\sqrt{\frac{T-p-h}{p}} \text{vec} \left(\hat{A}(p, h) - A(p, h) \right) \xrightarrow{d} N(0, \Sigma_\alpha)$ where $\Sigma_\alpha = \Gamma_p^{-1} \otimes \Omega_h$ and $\Omega_h = \Phi(I_h \otimes \Sigma_u)\Phi'$

where

$$\Phi = \begin{bmatrix} I_k & \mathbf{0} & \dots & \mathbf{0} \\ \Phi_1 & I_k & \dots & \mathbf{0} \\ \vdots & \vdots & \dots & \vdots \\ \Phi_{h-1} & \Phi_{h-2} & \dots & I_k \end{bmatrix}$$

(d) Let $\hat{\mathbf{u}}(p)_t \equiv y_t - \hat{\mathbf{m}} - \sum_{j=1}^p \hat{A}_j y_{t-j}$ so that $\hat{\Sigma}_u(p) = (T-p)^{-1} \sum_{t=1}^T \hat{\mathbf{u}}(p)_t \hat{\mathbf{u}}(p)_t'$ then

$\sqrt{T} \left(\hat{\Sigma}_u(p) - \Sigma_u \right) \rightarrow N(0, \Omega_\Sigma)$ where Ω_Σ is the covariance matrix of the residual covariance matrix.

Several results deserve comment. Technically speaking, condition (ii) in assumption 2 is required for asymptotic normality but not for consistency, where the weaker condition $p^2/T \rightarrow 0$, $T, p \rightarrow \infty$ is sufficient. Results (a)-(c) show that estimates of truncated models are consistent and asymptotically normal. Result (d) is useful if one prefers to rotate the vector of endogenous variables \mathbf{y}_t when providing structural interpretations for the forecast exercise. Here though, we abstain of such interpretation and provide the result only for completeness.

Next, denote with $\mathbf{y}_T(h)$ the forecast of the vector \mathbf{y}_{T+h} assuming the coefficients of the infinite order process (2) were known, that is

$$\mathbf{y}_T(h) = \mathbf{m} + \sum_{j=1}^{\infty} A_j \mathbf{y}_T(h-j)$$

where $\mathbf{y}_T(h-j) = \mathbf{y}_{T+h-j}$ for $h-j \leq 0$. Denote $\hat{\mathbf{y}}_T(h)$ the forecast that relies on coefficients estimated from a sample of size T and based on a finite order VAR or local projections, respectively

$$\begin{aligned} \hat{\mathbf{y}}_T(h) &= \hat{\mathbf{m}} + \sum_{j=1}^p \hat{A}_j \hat{\mathbf{y}}_T(h-j) \\ \hat{\mathbf{y}}_T(h) &= \hat{\mathbf{m}}_h + \sum_{j=0}^{p-1} \hat{A}_j^h \mathbf{y}_{T-j} \end{aligned}$$

were $\hat{\mathbf{y}}_T(h-j) = \mathbf{y}_{T+h-j}$ for $h-j \leq 0$. To economize in notation, we do not introduce a subscript that identifies how the forecast path was constructed as it should be obvious in the context of the derivations we provide. Then, define the forecast path for $h = 1, \dots, H$ by stacking each of the quantities $\hat{\mathbf{y}}_T(h)$, $\mathbf{y}_T(h)$, and \mathbf{y}_{T+h} as follows

$$\widehat{Y}_T(H)_{kH \times 1} = \begin{bmatrix} \widehat{\mathbf{y}}_T(1) \\ \vdots \\ \widehat{\mathbf{y}}_T(H) \end{bmatrix}; Y_T(H)_{kH \times 1} = \begin{bmatrix} \mathbf{y}_T(1) \\ \vdots \\ \mathbf{y}_T(H) \end{bmatrix}; Y_{T,H} = \begin{bmatrix} \mathbf{y}_{T+1} \\ \vdots \\ \mathbf{y}_{T+H} \end{bmatrix}.$$

Our interest is in finding the asymptotic distribution for $\widehat{Y}_T(H) - Y_{T,H} = [\widehat{Y}_T(H) - Y_T(H)] + [Y_T(H) - Y_{T,H}]$.

It should be clear that $[Y_T(H) - Y_{T,H}]$ does not depend on the estimation method and hence its mean-squared error can be easily verified to be

$$\Omega_H \equiv E[(Y_T(H) - Y_{T,H})(Y_T(H) - Y_{T,H})'] = \Phi(I_H \otimes \Sigma_u)\Phi'. \quad (5)$$

Furthermore, since parameter estimates are based on a sample of size T and hence \mathbf{u}_t for $t = p+h, \dots, T$ while the term $Y_T(H) - Y_{T,H}$ only involves \mathbf{u}_t for $T+1, \dots, T+H$, then it should be clear that to derive the asymptotic distribution of $[\widehat{Y}_T(H) - Y_T(H)]$, the asymptotic covariance of the forecast path will simply be the sum of the asymptotic covariance for this term and the mean-squared error in expression (5) and we do not have to worry about the covariance between these terms.

Corollary 1(a) and 1(b) and the observation that $\widehat{Y}_T(H)$ is simply a function of estimated parameters and predetermined variables is all we need to conclude that

$$\begin{aligned} & \sqrt{\frac{T-p-H}{p}} \text{vec}(\widehat{Y}_T(H) - Y_T(H)) \xrightarrow{d} N(0, \Psi_H) \\ \Psi_H & \equiv \frac{\partial \text{vec}(\widehat{Y}_T(H))}{\partial \text{vec}(\widehat{\mathbf{A}})} \Sigma_A \frac{\partial \text{vec}(\widehat{Y}_T(H))}{\partial \text{vec}(\widehat{\mathbf{A}})'} \end{aligned} \quad (6)$$

where Σ_A is the covariance matrix for $vec(\hat{\mathbf{A}})$; with $\hat{\mathbf{A}} = \hat{A}(p)$ for estimates from a VAR(p); and for estimates from local projections

$$\hat{\mathbf{A}} = \begin{bmatrix} \hat{A}(p, 1) \\ \vdots \\ \hat{A}(p, H) \end{bmatrix}. \quad (7)$$

We find it convenient to momentarily alter the order of our derivations and begin by examining forecasts from local projections first since these are linear functions of parameter estimates and hence can be obtained in a straight-forward manner.

First notice that $\hat{Y}_T(H) = \hat{\mathbf{A}}X_{T,p}$ and hence

$$\frac{\partial vec(\hat{Y}_T(H))}{\partial vec(\hat{\mathbf{A}})} = \frac{\partial vec(\hat{\mathbf{A}}X_{t,p})}{\partial vec(\hat{\mathbf{A}})} = \begin{pmatrix} X'_{T,p} \otimes I_{kH} \\ kH \times k^2Hp + kH \end{pmatrix} \quad (8)$$

which combined with corollary 1(c) results in

$$\sqrt{\frac{T-p-H}{p}} \left(vec(\hat{\mathbf{A}} - \mathbf{A}) \right) \xrightarrow{d} N(0, \Sigma_A) \quad (9)$$

$$\begin{matrix} \Sigma_A \\ k^2Hp + kH \times k^2Hp + kH \end{matrix} = \begin{matrix} \Gamma_p^{-1} \otimes \Omega_H; \\ \Omega_H \\ kH \times kH \end{matrix} = \mathbf{\Phi} (I_H \otimes \Sigma_u) \mathbf{\Phi}'$$

Putting together expressions (6), (5), (8) and (9), we arrive at the following corollary.

Corollary 2 *Under assumptions 1 and 2 and expressions (6), (5), (8) and (9), the asymptotic distribution of the forecast path generated with the local projections approach described*

in assumption 1 is

$$\begin{aligned}
& \sqrt{\frac{T-p-H}{p}} \text{vec} \left(\hat{Y}_T(H) - Y_{T,H} \right) \xrightarrow{d} N(\mathbf{0}; \Xi_H) \\
\Xi_H &= \left\{ \frac{p}{T-p-H} \Omega_H + \Psi_H \right\} \\
\Omega_H &= \Phi(I_H \otimes \Sigma_u) \Phi' \\
\Psi_H &= (X'_{T,p} \otimes I_{kH}) [\Gamma_p^{-1} \otimes \Omega_H] (X_{T,p} \otimes I_{kH})
\end{aligned} \tag{10}$$

In practice, all population moments can be substituted by their conventional sample counterparts.

We now return to the more involved derivation of the asymptotic distribution of the forecast path when the forecasts are generated by the VAR(p) in expression (3). For this purpose, we find it easier to work with each element of the vector $\hat{Y}_T(H)$ individually, so that we begin by examining the derivation of

$$\begin{aligned}
& \sqrt{\frac{T-p-H}{p}} \text{vec} (\hat{\mathbf{y}}_T(h) - \mathbf{y}_T(h)) \xrightarrow{d} N(\mathbf{0}; \Psi_{h,h}) \\
\Psi_{h,h} &= \frac{\partial \text{vec} (\hat{\mathbf{y}}_T(h))}{\partial \text{vec} (\hat{A}(p))} \Sigma_a \frac{\partial \text{vec} (\hat{\mathbf{y}}_T(h))}{\partial \text{vec} (\hat{A}(p))}
\end{aligned}$$

where we remind the reader that from corollary 1(b), $\Sigma_a = \Gamma_p^{-1} \otimes \Sigma_u$. In general, notice that

$$\Psi_{i,j} = \frac{\partial \text{vec} (\hat{\mathbf{y}}_T(i))}{\partial \text{vec} (\hat{A}(p))} \Sigma_a \frac{\partial \text{vec} (\hat{\mathbf{y}}_T(j))}{\partial \text{vec} (\hat{A}(p))}$$

which is all we need to construct all the elements in the asymptotic covariance matrix of $\hat{Y}_T(H)$, namely Ψ_H . An expression for $\hat{\mathbf{y}}_T(h)$ generated from the VAR(p) in expression (3) can be obtained as

$$\widehat{\mathbf{y}}_T(h) = JB^h X_{T,p}$$

where B simply stacks the $\text{VAR}(p)$ coefficients in companion form and J is a selector matrix, both of which are

$$\begin{aligned} B_{kp+1 \times kp+1} &= \begin{pmatrix} 1 & 0 & 0 & \dots & 0 & 0 \\ \mathbf{m} & A_1 & A_2 & \dots & A_{p-1} & A_p \\ 0 & I_k & 0 & \dots & 0 & 0 \\ 0 & 0 & I_k & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & I_k & 0 \end{pmatrix}, \\ J_{k \times kp+1} &= (\mathbf{0}_{k \times 1}, I_k, \mathbf{0}_k, \dots, \mathbf{0}_k)_{k \times k}. \end{aligned}$$

Therefore, notice that

$$\frac{\partial \text{vec}(\widehat{\mathbf{y}}_T(h))}{\partial \text{vec}(\widehat{A}(p))} = \frac{\partial \text{vec}(JB^h X_{t,p})}{\partial \text{vec}(\widehat{A}(p))} = \sum_{i=0}^{h-1} X'_{T,p} (B')^{h-1-i} \otimes \Pi_i, \quad \Pi_i = JB^i J'.$$

The following corollary characterizes the asymptotic distribution of $\text{VAR}(p)$ generated forecasts paths.

Corollary 3 *Under assumptions 1 and 2, the asymptotic distribution of the forecast path*

$\widehat{Y}_T(H)$ generated from the $VAR(p)$ in expression (3) is given by

$$\begin{aligned}
& \sqrt{\frac{T-p-H}{p}} \text{vec} \left(\widehat{Y}_T(H) - Y_{T,H} \right) \xrightarrow{d} N(\mathbf{0}; \Xi_H) \\
\Xi_H &= \left\{ \frac{p}{T-p-H} \Omega_H + \Psi_H \right\} \\
\Omega_H &= \Phi(I_H \otimes \Sigma_u) \Phi' \\
\Psi_{i,j} &= \frac{p}{T-p-H} \sum_{k=0}^{i-1} \sum_{s=0}^{j-1} E(X'_{T,p} (B')^{i-1-k} \Gamma_p^{-1} B^{j-1-s} X_{T,p}) \otimes \Pi_k \Sigma_u \Pi'_s \\
&= \frac{p}{T-p-H} \sum_{k=0}^{i-1} \sum_{s=0}^{j-1} \text{tr}((B')^{i-1-k} \Gamma_p^{-1} B^{j-1-s} \Gamma_p) \Pi_k \Sigma_u \Pi'_s
\end{aligned} \tag{11}$$

In practice all moment matrices can be substituted by their sample counterparts as usual.

Thus, corollaries 2 and 3 provide the necessary results on the joint predictive density of the path forecasts. The next sections exploit these results to introduce new methods of path forecast comparison and predictive ability testing.

3 Model Comparison: Mean Squared Forecast Path

The most commonly used metric of forecast model comparison is the mean squared forecast error (*MSFE*). Given a sample of $1, \dots, T$ observations available for estimation and $T + 1, \dots, T + N$ observations available for forecast evaluation, this metric is constructed as

$$MSFE_h = \frac{1}{N} \sum_{i=1}^N (\widehat{\mathbf{y}}_{T+i}(h) - \mathbf{y}_{T+i+h})' (\widehat{\mathbf{y}}_{T+i}(h) - \mathbf{y}_{T+i+h})$$

so that an absolute comparison of the overall predictive merits between two competing models for a particular forecast horizon h can be directly obtained by comparing their respective $MSFE_h$.

However, often times a model that predicts well at short horizons will predict badly a longer horizons (and vice versa), thus making an assessment of overall predictive performance difficult. A natural way to overcome this difficulty is to construct a metric that evaluates entire forecast paths jointly. Clements and Hendry (1993) suggest to base comparison on the determinant of the forecast error second moment matrix pooled across horizons of interest, which they call *generalized forecast error second moment* (*GFESM*). Compared to the standard *MSFE*, the *GFESM* has the advantage of being invariant to non-singular, scale-preserving linear transformations. However, *GFESM* does not provide a natural basis on which to build tests of relative predictive ability.

Instead, suppose that the $1, \dots, H$ forecast path $\hat{Y}_T(H)$ has an asymptotic distribution given by

$$\sqrt{T} \left(\hat{Y}_T(H) - Y_{T,H} \right) \xrightarrow{d} N(0; \Xi_H).$$

Examples of such a result are corollaries 2 and 3. We know then that the associated Wald statistic

$$W_H = T \left(\hat{Y}_T(H) - Y_{T,H} \right)' \Xi_H^{-1} \left(\hat{Y}_T(H) - Y_{T,H} \right) \xrightarrow{d} \chi_{kH}^2$$

provides a natural metric of distance between $\hat{Y}_T(H)$ and $Y_{T,H}$. This metric operates at two levels: (1) the relative efficiency with which each forecast is generated; and (2) the degree of correlation between forecasts. Based on this distance metric, the measure in equivalent units to the *MSFE_h* is

$$MSFP_{1,H} = \left(\frac{1}{HN} \sum_{i=1}^N \left(\hat{Y}_{T+i}(H) - Y_{T+i,H} \right)' \hat{\Lambda}_H^{-1} \left(\hat{Y}_{T+i}(H) - Y_{T+i,H} \right) \right)$$

where Λ_H is the correlation matrix associated to Ξ_H and $\hat{\Lambda}_H \xrightarrow{p} \Lambda_H$ if $\hat{\Xi}_H \xrightarrow{p} \Xi_H$ such as in corollaries 2 and 3; and $MSFP_{1,H}$ stands for the mean squared forecast path over horizons $1, \dots, H$.

For $H = 1$, then $MSFP_{1,1} = MSFE_1$. Similarly, if forecasts at horizons $1, \dots, H$ are uncorrelated (Ξ_H is diagonal) then $MSFP_{1,H} = \frac{1}{H} \sum_{h=1}^H MSFE_h$, that is, an average of the MSFE over $1, \dots, H$. To get a better sense, suppose the data were generated by the simple AR(1) model

$$y_t = \rho y_{t-1} + \varepsilon_t \quad \varepsilon_t \sim N(0, \sigma^2)$$

and for simplicity assume that ρ is known rather than estimated. Next, consider the forecast path for $h = 1, 2$

$$\hat{Y}_T(2) = \begin{bmatrix} \rho y_T \\ \rho^2 y_T \end{bmatrix}; Y_{T,2} = \begin{bmatrix} \rho y_T + \varepsilon_{T+1} \\ \rho^2 y_T + \varepsilon_{T+2} + \rho \varepsilon_{T+1} \end{bmatrix}$$

then clearly

$$\hat{Y}_T(2) - Y_{T,2} = \begin{bmatrix} \varepsilon_{T+1} \\ \varepsilon_{T+2} + \rho \varepsilon_{T+1} \end{bmatrix} \sim N \left(0; \sigma^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 + \rho^2 \end{bmatrix} \right)$$

The $MSFP_{1,H}$ for a predictive sample of N observations is then

$$\begin{aligned}
MSFP_{1,2}^{AR} &= \frac{1}{2N} \sum_{i=1}^N \left(\widehat{Y}_{T+i}(2) - Y_{T+i,2} \right)' \begin{bmatrix} 1 & \rho \\ \rho & 1 + \rho^2 \end{bmatrix}^{-1} \left(\widehat{Y}_{T+i}(2) - Y_{T+i,2} \right) \\
&= \frac{1}{2N} \sum_{i=1}^N \begin{pmatrix} \varepsilon_{T+i+1} & \varepsilon_{T+i+2} + \rho \varepsilon_{T+i+1} \end{pmatrix} \begin{bmatrix} 1 & \rho \\ \rho & 1 + \rho^2 \end{bmatrix}^{-1} \begin{pmatrix} \varepsilon_{T+i+1} \\ \varepsilon_{T+i+2} + \rho \varepsilon_{T+i+1} \end{pmatrix} \\
&= \frac{1}{2N} \sum_{i=1}^N (\varepsilon_{T+i+1}^2 + \varepsilon_{T+i+2}^2) \xrightarrow{p} \sigma^2 \text{ as } N \rightarrow \infty
\end{aligned}$$

That is, absent parameter estimation uncertainty and given that ε_t are *i.i.d.* normal, then the $MSFP_{1,2}^{AR}$ is simply the predictive sample variance of the shocks hitting the model at each forecast horizon. In other words, this is a pure measure of how well the model fits at each horizon, *distilled from the correlation between how the forecasts are constructed over time.*¹

The next section exploits the Wald metric on which $MSFP$ is based to extend available tests of relative predictive ability between models.

4 Tests of Path Predictive Ability

The $MSFP_{1,H}$ metric allows one to easily determine which of two competing forecasting models performs best, *in the available sample of data*. If we want to determine whether differences in forecasting performance are statistically significant, we need to recast common tests of predictive ability in terms of paths. The literature on comparing the predictive ability of competing forecasts given general loss functions was initiated by Diebold and Mariano (1995) and further formalized in West (1996); McCracken (2000); Clark and McCracken

¹ Parenthetically, notice that when one omits parameter estimation uncertainty, there is no difference between forecasts made by iterating the AR(1) specification or by local projections so that this result is not method-dependent.

(2001); Corradi, Swanson and Olivetti (2001); and Chao, Corradi and Swanson (2001), among others. Giacomini and White (2006) extend this literature even further by considering conditional predictive ability tests based on examining null hypotheses that are expressed in terms of conditional expected forecast loss functions rather than unconditionally, as had previously been done. The literature is obviously very extensive so we cannot presume to explore the details of every possible contingency when testing path predictive ability. However, we think the underlying principle can be succinctly presented for the most common testing scenario, and we leave for further research appropriate generalizations.

Accordingly, suppose interest is in comparing the accuracy of two competing forecasting models whose forecasted paths are $\widehat{Y}_T^j(H)$ for $j = 1, 2$ indicating the method. We can then think of the test of equal predictive ability as a Hausman test, where the null hypothesis we are interested in testing is

$$H_0 : p \lim vec \left[\overline{\widehat{Y}}_T^1(H) - \overline{\widehat{Y}}_T^2(H) \right] = 0, \quad N \rightarrow \infty$$

where T is the sample size available for estimation, N is the sample size available for forecast evaluation and

$$\overline{\widehat{Y}}_T^j(H) = \frac{1}{N} \sum_{i=1}^N \widehat{Y}_{T+i}^j(H); \quad j = 1, 2.$$

The key element to derive the asymptotic properties of the test is then to establish sufficient conditions so that an appropriate central limit theorem holds for

$$\sqrt{N} vec \left[\overline{\widehat{Y}}_T^1(H) - \overline{\widehat{Y}}_T^2(H) \right] \xrightarrow{d} N(0; \Upsilon_H) \quad (12)$$

and then the null hypothesis can be evaluated with the statistic

$$PDMW_H = N \left[\widehat{Y}_T(H) - \overline{Y}_T(H) \right]' \Upsilon_H^{-1} \left[\widehat{Y}_T(H) - \overline{Y}_T(H) \right] \rightarrow \chi_{kH}^2.$$

Several remarks are worth making. First, notice that if $T, N \rightarrow \infty$ (as is done in West, 1996) then parameter estimation uncertainty vanishes. As an example, consider the two forecasting approaches and DGP we have considered in previous sections. Generically speaking, expressions (10) and (11) have the general format

$$\begin{aligned} \sqrt{T} \left(\widehat{Y}_T^j(H) - Y_{T,H} \right) &\xrightarrow{d} N \left(0, \Xi_H^j \right); \\ \Xi_H^j &= \left[\frac{1}{T} \Omega_H + \Psi_H^j \right] \end{aligned} \tag{13}$$

so that as $T \rightarrow \infty$, $\Psi_H^j \xrightarrow{p} 0$ and expression (12) is a natural consequence of (13). Hence, assumptions that guarantee the result (13), such as assumptions 1 and 2 in our paper, are sufficient to guarantee (12). However, when T is fixed (as in Giacomini and White, 2006) the asymptotic distribution of the test statistics will reflect non-vanishing estimation uncertainty. In that case, what is needed are conditions that guarantee that $\text{vec} \left[\widehat{Y}_T(H) - \overline{Y}_T(H) \right]$ is a martingale difference sequence with appropriate mixing conditions that ensure a central limit theorem is available to derive expression (12). Often times the literature has opted to make these assumptions primitive with respect to the asymptotic framework and we refer the reader to Giacomini and White (2006) for very general conditions that ensure the result holds. Second, notice that $\Upsilon_H = \Xi_H^1 + \Xi_H^2 - \Xi_H^{1,2} - \Xi_H^{2,1}$ where $\Xi_H^{i,j}$ denotes the covariance between the forecast paths obtained by the two methods being considered. Under assumptions 1 and 2, as $T \rightarrow \infty$ even with N fixed, it would be possible to obtain consistent estimates of Ξ_H^1 and Ξ_H^2 with the formulas that we provide in corollaries 2 and 3 but in general, it is often

not possible to obtain closed-form expressions for $\Xi_H^{i,j}$. Under assumptions 1 and 2, then a consistent estimate of Υ_H (as $N \rightarrow \infty$) is

$$\begin{aligned}\hat{\Upsilon}_H &= \frac{1}{N} \sum_{i=1}^N \left(\hat{Y}_T^1(H) - \bar{Y}_T^1(H) \right) \left(\hat{Y}_T^1(H) - \bar{Y}_T^1(H) \right)' \\ &\quad + \frac{1}{N} \sum_{i=1}^N \left(\hat{Y}_T^2(H) - \bar{Y}_T^2(H) \right) \left(\hat{Y}_T^2(H) - \bar{Y}_T^2(H) \right)' \\ &\quad - \frac{1}{N} \sum_{i=1}^N \left(\hat{Y}_T^1(H) - \bar{Y}_T^1(H) \right) \left(\hat{Y}_T^2(H) - \bar{Y}_T^2(H) \right)' \\ &\quad - \frac{1}{N} \sum_{i=1}^N \left(\hat{Y}_T^2(H) - \bar{Y}_T^2(H) \right) \left(\hat{Y}_T^1(H) - \bar{Y}_T^1(H) \right)'\end{aligned}$$

Assumptions that allow for heterogeneity of $\hat{Y}_{T+i}^j(H)$ across the evaluation sample $i = T + 1, \dots, T + N$ will instead require heteroskedasticity and autocorrelation (HAC) robust versions of $\hat{\Upsilon}_H$, such as a Newey-West covariance matrix estimator.

5 Simultaneous Confidence Regions for Forecast Paths

This section considers the problem of constructing a simultaneous confidence region for the forecast path of the j^{th} variable in the k -dimensional system we have so far examined. Under assumptions 1 and 2, corollaries 2 and 3 show that the asymptotic distribution of $\hat{Y}_T(H)$ is (with the obvious simplifications relative to (10) or (11)):

$$\sqrt{T} \left(\hat{Y}_T(H) - Y_{T,H} \right) \xrightarrow{d} N(\mathbf{0}; \Xi_H). \quad (14)$$

Let $S_j \equiv (I_H \otimes \mathbf{e}_j)$ where \mathbf{e}_j is a $1 \times k$ vector of zeros with a 1 in the j^{th} column. Then the asymptotic distribution for the forecast path of the j^{th} variable in (14) is readily seen to be

$$\sqrt{T} \left(\widehat{Y}_{j,T}(H) - Y_{j,T,H} \right) \xrightarrow{d} N \left(\mathbf{0}; \boldsymbol{\Xi}_{j,H} \right), \quad (15)$$

where $\widehat{Y}_{j,T}(H) = S_j \widehat{Y}_T(H)$; $Y_{j,T,H} = S_j Y_{T,H}$; and $\boldsymbol{\Xi}_{j,H} = S_j \boldsymbol{\Xi}_H S_j'$. The derivations we are about to present do not depend on how one arrives at expression (15). Therefore, to make the results more general, we take expression (15) as our primitive assumption so as to accommodate forecasting environments other than those implied by assumptions 1 and 2.

The conventional approach to reporting forecasting uncertainty consists of displaying two standard-error bands constructed from the square roots of the diagonal entries of $\boldsymbol{\Xi}_{j,H}$. The confidence region described by these bands is therefore equivalent to testing a joint null hypothesis with the collection of t-statistics associated to the individual elements of the joint null. It is easy to see that such an approach ignores the simultaneous nature of the problem and any correlation that may exist among the forecasts across horizons, thus providing incorrect probability coverage.

Under expression (15), the Wald principle suggests that a joint null hypothesis on $Y_{j,T,H}$ can be tested with the statistic

$$W_H = T \left(\widehat{Y}_{j,T}(H) - Y_{j,T,H} \right)' \boldsymbol{\Xi}_{j,H}^{-1} \left(\widehat{Y}_{j,T}(H) - Y_{j,T,H} \right) \xrightarrow{d} \chi_H^2 \quad (16)$$

so that a confidence region at an α significance level is represented by the values of $Y_{j,T,H}$ that satisfy

$$P \left[W_H \leq c_\alpha^2 \right] = 1 - \alpha$$

where c_α^2 is the critical value of a random variable distributed χ_H^2 at a $1 - \alpha$ confidence level.

This confidence region is a multi-dimensional ellipsoid that in general, is too complicated to display graphically and makes communication of forecast uncertainty difficult. However, for $H = 2$, this region can be displayed in two-dimensional space as is done in figure 1.

The top panel of figure 1 displays the 95% confidence region associated to a one- and two-period ahead forecasts from an AR(1) model with known autoregressive coefficient $\rho = 0.75$ and $\sigma = 1$. Overlaid on this ellipse is the traditional two standard-error box. The figure makes clear why this box provides inappropriate probability coverage: it contains/excludes forecast paths with less/more than 5% chance of being observed.

In order to reconcile the inherent difficulty of displaying multi-dimensional ellipsoids with the inadequate probability coverage provided by the more easily displayed marginal error bands, we construct a simultaneous rectangular region with Scheffé's (1953) S-method of simultaneous inference (see also Lehmann and Romano, 2005). Briefly, the intuition of the method is to exploit the Cauchy-Schwarz inequality to transform the Wald statistic from L_2 -metric into L_1 -metric to facilitate construction of a rectangular confidence interval.

Notice that the covariance matrix of $\hat{Y}_{j,T}(H)$ is positive-definite and symmetric and hence admits a Cholesky decomposition $T^{-1}\Xi_{j,H} = PP'$, where P is a lower triangular matrix. The passage of time provides a natural and unique ordering principle so that P is obtained unambiguously. Notice then that

$$\begin{aligned}
P \left[T \left(\widehat{Y}_{j,T}(H) - Y_{j,T,H} \right)' \Xi_{j,H} \left(\widehat{Y}_{j,T}(H) - Y_{j,T,H} \right) \leq c_\alpha^2 \right] &= 1 - \alpha \\
P \left[\left(\widehat{Y}_{j,T}(H) - Y_{j,T,H} \right)' (PP')^{-1} \left(\widehat{Y}_{j,T}(H) - Y_{j,T,H} \right) \leq c_\alpha^2 \right] &= 1 - \alpha \\
P \left[\widehat{Z}_{j,T}(H)' \widehat{Z}_{j,T}(H) \leq c_\alpha^2 \right] &= 1 - \alpha \\
P \left[\sum_{h=1}^H \widehat{z}_{j,T}(h)^2 \leq c_\alpha^2 \right] &= 1 - \alpha \quad (17)
\end{aligned}$$

where $\widehat{Z}_{j,T}(H) = P^{-1} \widehat{Y}_{j,T}(H)$ and $\widehat{z}_{j,T}(h) \rightarrow N(0, 1)$ and independent across h . Consider now the problem of constructing the rectangular confidence region for the average path-average forecast

$$P \left[\left| \sum_{h=1}^H \frac{\widehat{z}_{j,T}(h)}{h} \right| \leq \delta_\alpha \right] = 1 - \alpha.$$

A direct consequence of Bowden's (1970) lemma is that

$$\max \left\{ \frac{\left| \sum_{h=1}^H \frac{\widehat{z}_{j,T}(h)}{h} \right|}{\sqrt{\sum_{h=1}^H \frac{1}{h^2}}} : |h| < \infty \right\} = \sqrt{\sum_{h=1}^H \widehat{z}_{j,T}(h)^2}$$

which can be applied directly to the bottom line of expression (17) to obtain

$$P \left[\left| \sum_{h=1}^H \frac{\widehat{z}_{j,T}(h)}{h} \right| \leq \sqrt{\frac{c_\alpha^2}{H}} \right] = 1 - \alpha, \quad (18)$$

which in turn implies that $\delta_\alpha = \sqrt{\frac{c_\alpha^2}{H}}$. Expression (18) and $\widehat{Z}_{j,T}(H) = P^{-1} \widehat{Y}_{j,T}(H)$ imply that a simultaneous confidence region for the forecast path $\widehat{Y}_{j,T}(H)$ can then be constructed as

$$\widehat{Y}_{j,T}(H) \pm \delta_\alpha P \mathbf{i}_H \quad (19)$$

where \mathbf{i}_H is an $H \times 1$ vector of ones. We call these bands Scheffé confidence bands to distinguish them from the usual two (marginal) standard error bands commonly reported.

One way to gain intuition about the Scheffé bands is to establish their relation to traditional marginal error bands. In a traditional error band, the boundaries of the band represent a shift from the mean of the distribution (the parameter estimate) in proportion to its variance. Thus, the boundary is the appropriately scaled $1 - \alpha$ critical value of the standard normal density, that is for example, $\hat{y}_{j,T}(h) \pm z_{\alpha/2} \hat{\Xi}_{j,(h,h)}^{1/2}$. Similarly, consider now a *simultaneous* shift in all the elements of the forecast path in proportion to their variances. What would the appropriate critical value be? It is easier to answer this question with the orthogonal coordinate system $\hat{Z}_{j,T}(H)$. From expression (17) and denoting this critical value δ_α , then δ_α must meet the condition

$$P [\delta_\alpha^2 + \frac{H}{H} + \delta_\alpha^2 = c_\alpha^2] = 1 - \alpha$$

which implies that $\delta_\alpha = \sqrt{\frac{c_\alpha^2}{H}}$. In two dimensions, both panels of figure 1 display diagonal intersecting the origin of both ellipses. The slopes of these diagonals reflect the relative variance of each forecast (in the bottom panel the normalization ensures the variances are the same so the diagonal is the 45 degree line) and represent the $\pm\delta_\alpha$ for all values of α . The Cholesky factor P therefore provides the appropriate scaling for δ_α since it not only scales the orthogonal system by the individual variances of its elements but also accounts for their correlation.

Several results deserve comment. First, when $H = 1$ so that we are considering a one-period ahead forecast, then $c_{0.05}^2 = 3.84$ for a χ_1^2 random variable and hence $\delta_\alpha = \sqrt{3.84} =$

1.96. In this case, $P = \sigma_1$ so that the rectangular confidence interval obtain by Scheffé's S-method corresponds to the traditional two standard-error band. However, when $H = 2$, then $\delta_\alpha = 1.73$, not the usual 1.96. Second, because the relation between the L_2 -norm implied by the Wald statistic and the rectangular region implied by the associated L_1 -norm holds by Hölder's inequality (rather than with equality), the probability coverage is more conservative. Third, an alternative approach is to construct confidence intervals with Bonferroni's inequality. This inequality suggests that a $(1 - \frac{\alpha}{H})$ confidence interval for $y_{j,T}(h)$, $h = 1, \dots, H$, then the union of these confidence intervals generates a region that includes $Y_{j,T,H}$ with at least $(1 - \alpha)$ probability. Specifically, the Bonferroni confidence region (BCR) is

$$\hat{Y}_{j,T}(H) \pm z_{\alpha/2H} \times \text{diag}(\Xi_{j,H}),$$

where $z_{\alpha/2H}$ denotes the critical value of a standard normal random variable at a $\alpha/2H$ significance level. Thus, the BCR can be significantly more conservative than our simultaneous confidence region, specially when the correlation between forecasts across horizons is low.

In addition, we offer two complementary ways to report uncertainty about the forecast. The first is to notice that it is easy to construct a fan chart with the quantiles of the joint predictive density simply by calculating the simultaneous rectangular regions associated with the values that c_α^2 and δ_α take for different values of α . An example of such a chart is provided below in the empirical section. The second measure is based on the following observation. Notice that $\Xi_{j,H} = PP' = QDQ'$ where Q is lower triangular with ones in the main diagonal and D is a diagonal matrix. Hence, the Wald statistic in expression (16) can be rewritten as

$$\begin{aligned}
W_H &= T \left(\widehat{Y}_{j,T}(H) - Y_{j,T,H} \right)' \Xi_{j,H}^{-1} \left(\widehat{Y}_{j,T}(H) - Y_{j,T,H} \right) \\
&= \left(\widehat{Y}_{j,T}(H) - Y_{j,T,H} \right)' (QDQ')^{-1} \left(\widehat{Y}_{j,T}(H) - Y_{j,T,H} \right) \\
&= \widetilde{Z}_{j,T}(H)' D^{-1} \widetilde{Z}_{j,T}(H) \\
&= \sum_{h=1}^H \frac{\widetilde{z}_{j,T}(h)^2}{d_{hh}} = \sum_{h=1}^H t_{h|h-1, \dots, 1}^2 \rightarrow \chi_H^2
\end{aligned}$$

where $\widetilde{Z}_{j,T}(H) = Q^{-1} \left(\widehat{Y}_{j,T}(H) - Y_{j,T,H} \right)$ is the unstandardized version of $\widehat{Z}_{j,T}(H)$; and d_{hh} is the h^{th} diagonal entry of D , which is the variance of $\widetilde{z}_{j,T}(h)$. In other words, the Wald statistic W_H of the joint null on $Y_{j,T,H}$ is equivalent to the sum of the squares of the conditional t -statistics of the individual nulls of significance of the forecast path at time h given the path from 1 to $h-1$ for $h = 1, \dots, H$. An implication of this result is that an $1 - \alpha$ confidence region that sterilizes the uncertainty about the forecast path up to time $h-1$ and summarizes the uncertainty about the h horizon forecast alone, can be easily constructed with the bands

$$\widehat{Y}_{j,T}(H) \pm z_{\alpha/2} \times \text{diag}(D)$$

where $z_{\alpha/2}$ refers to the critical value of a standard normal random variable at a $\alpha/2$ significance level.

A simple example provides further intuition about the relation between the different confidence regions discussed in this section. Suppose the data are generated by the AR(1) model previously discussed in section 3, then the 95% confidence ellipse results from the associated Wald statistic in expression (16), that is

$$W_2 = \frac{1}{\sigma^2} \left(\widehat{Y}_T(2) - Y_{T,2} \right)' \begin{bmatrix} 1 & \rho \\ \rho & 1 + \rho^2 \end{bmatrix}^{-1} \left(\widehat{Y}_T(2) - Y_{T,2} \right) \leq 6$$

since $c_{0.05}^2 \simeq 6$ for a χ_2^2 random variable. This is the ellipse displayed in the top panel of figure 1 for $\sigma = 1$ and $\rho = 0.75$. For this example, the traditional two standard-error box is given by $[-1.96, 1.96]$ and $[-2.45, 2.45]$.

The Cholesky decomposition of this forecast path's covariance matrix is

$$\sigma^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 + \rho^2 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ \rho & 1 \end{bmatrix} \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix} \begin{bmatrix} 1 & \rho \\ 0 & 1 \end{bmatrix} \quad (20)$$

and it is clear that the orthogonal path's covariance matrix is

$$\begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}$$

that is, this covariance simply summarizes the fact that the only source of forecast uncertainty in this simple example comes from the *i.i.d.* shock, ε_t (to see it more clearly suppose that $\rho = 0$). The 95% confidence circle associated to the orthogonalized forecast path is displayed in the bottom panel of figure 1. The associated rectangular region can be easily constructed by noticing that $\delta_\alpha = \sqrt{c_\alpha^2/2} = 1.73$ and hence the box is given by $[-1.73, 1.73]$ and $[-1.73, 1.73]$. Parenthetically, these also correspond to the conditional two standard error bands.

In order to obtain the simultaneous, rectangular 95% confidence region we need to translate the rectangular box in the orthogonal coordinate system back to the original coordinate

system of the forecast path. The Cholesky factor when $\rho = 0.75$ can be easily obtained from expression (20), which multiplied by the orthogonal rectangular values

$$\pm \begin{bmatrix} 1 & 0 \\ 0.75 & 1 \end{bmatrix} \begin{bmatrix} 1.73 \\ 1.73 \end{bmatrix} = \pm \begin{bmatrix} 1.73 \\ 3.03 \end{bmatrix}$$

delivers the simultaneous 95% rectangular region $[-1.73, 1.73]$ and $[-3.03, 3.03]$. Notice that compared to the traditional two standard-error box $[-1.96, 1.96]$ and $[-2.45, 2.45]$, Scheffé's (1953) method produces a confidence region in which the first period's forecast uncertainty is narrower but the second period's is wider. Figure 2 translates all of these bands discussed into a more traditional format to make the interpretation more transparent.

6 Conditional Path Forecasts

This section discusses a final method of path forecast evaluation that exploits the asymptotic Gaussian approximation of the joint predictive density in corollaries 2 and 3; and properties of the multivariate normal distribution and linear projections. The problem that we have in mind is that of constructing forecast paths conditional on alternative hypothetical paths for a subset of variables in the k -dimensional system being considered. An example perhaps provides better intuition about what we mean.

Suppose a policy maker is confronted with a set of path forecasts $\hat{Y}_T(H)$ about the future behavior of output growth, inflation, oil prices, exchange rates, interest rates, and so on. Given these forecasts, suppose the policy maker wants to stress the model and examine how would these macroeconomic forecasts vary if, for example, a path of oil prices different than that predicted were to take place. It turns out that when the joint predictive density is

Gaussian, not only is it simple to obtain what the conditional forecast paths would be, it is straightforward to calculate the associated conditional predictive density. Waggoner and Zha (1999) develop Bayesian methods to compute this distribution in finite samples for VARs whereas Leeper and Zha (2003) further investigate projections based on hypothetical paths of monetary policy.

It is important to remark that this section does not specifically address parameter stability in the face of variability in the hypothetical paths in the sense discussed by Lucas (1976). Such questions are reserved for future research. Rather, we examine hypothetical paths drawn from the joint predictive density and therefore, a natural starting point is to examine the likelihood of the hypothetical paths proposed. For this reason, we define the some additional notation.

Specifically, let $k = k_0 + k_1$ where k_0 is the dimension of the subvector of variables whose conditional path forecasts we want to calculate, and k_1 is the subvector of variables whose hypothetical paths are the conditioning set. Define the selector matrices $S_0 = (I_H \otimes E_0)$ and $S_1 = (I_H \otimes E_1)$ where E_0 is a $k_0 \times k$ matrix formed with the k_0 rows of I_k associated with the variables whose conditional paths we wish to calculate, and E_1 is a $k_1 \times k$ matrix whose rows are the k_1 rows of I_k associated with the variables whose hypothetical paths provide the conditioning set. Therefore, if $\hat{Y}_T^c(H)$ denotes the $kH \times 1$ vector of conditional forecasts, $\hat{Y}_T^0(H) = S_0 \hat{Y}_T^c(H)$ is the $k_0 H \times 1$ vector of forecasts conditional on the hypothetical paths $Y_T^1(H) = S_1 \hat{Y}_T^c(H)$, which are of dimension $k_1 H \times 1$ and where the “hat” is omitted because the hypothetical paths are not estimated but rather are assumed.

Recall that under assumptions 1 and 2, corollaries 2 and 3 suggest that the asymptotic

predictive density of $\widehat{Y}_T(H)$ is

$$\sqrt{T} \left(\widehat{Y}_T(H) - Y_{T,H} \right) \xrightarrow{d} N(\mathbf{0}; \Xi_H)$$

although we remark that different assumptions and forecast environments could produce a similar asymptotic result. The likelihood of the hypothetical paths $Y_T^1(H)$ can be evaluated by testing the null hypothesis $H_0 : S_1 Y_{T,H} = Y_T^1(H)$ with the Wald statistic

$$W_1^c = T \left(S_1 \widehat{Y}_T(H) - Y_T^1(H) \right)' (S_1 \Xi_H S_1')^{-1} \left(S_1 \widehat{Y}_T(H) - Y_T^1(H) \right) \xrightarrow{d} \chi_{k_1 H}^2.$$

Hence, one minus the p-value associated with W_1^c is a measure of the distance, in probability units, between the predicted paths and the hypothetical paths of the k_1 variables. Low p -values (i.e. toward the direction of rejecting the null) stress the conditional forecasting exercise toward paths for the k_1 variables that have been rarely observed in the historical sample. In such situations, the resulting conditional forecasts are likely to be more problematic because we are asking about areas where the model has not been trained by the sample.

Once the likelihood of the hypothetical values has been assessed, we want to calculate the conditional forecast paths $\widehat{Y}_T^0(H)$ and their predictive distribution. Standard properties of linear projections and the multivariate Gaussian distribution are all is needed to conclude that

$$\widehat{Y}_T^0(H) = \widehat{Y}_T(H) + S_0 \Xi_H S_1' (S_1 \Xi_H S_1')^{-1} (Y_T^1(H) - S_1 \widehat{Y}_T(H))$$

with a Gaussian predictive density whose covariance matrix is

$$\Xi_H^0 = S_0 \Xi_H S_0' - S_0 \Xi_H S_1' (S_1 \Xi_H S_1')^{-1} S_1 \Xi_H S_0.$$

It is worth remarking that the Gaussian approximation and Ξ_H^0 is all that is needed to construct simultaneous rectangular confidence regions, conditional error bands and Wald tests of joint hypotheses with the techniques described in previous sections. The next section illustrates all of these techniques with an empirical illustration.

7 Empirical Illustration

This section contains two illustration of the techniques presented above. The first application serves to illustrate the newly introduced measures of forecast path comparison and examines the role of monetary aggregates in forecasting inflation, a topic of much recent debate in central bank circles. The second application examines path forecasts for a system of U.S. macroeconomic variables and provides an example of counterfactual simulation.

7.1 Do Monetary Aggregates Help Forecast Inflation?

In order to achieve monetary policy goals, central banks constantly monitor risks to price stability. Since its inception, the European Central Bank (ECB) has used a so-called “two pillar” approach that gives a specific role to monetary aggregates over other economic indicators when forecasting inflation. Naturally, the question we ask in this section is whether monetary aggregates help forecast U.S. inflation since the U.S. Federal Reserve (Fed) does not reserve a special role for monetary aggregates over other predictors.

Our objective is very modest: to evaluate whether forecasts of core consumer price inflation (measured by the consumer price index less food and energy) improve when either the

M2 or the M3 money stock growth measures are added alongside the industrial production index growth rate and the federal funds rate – inflation, output and interest rates being the constituent variables of practically any macroeconomic model of the economy. All variables are observed monthly and are seasonally adjusted from January 1985 to January 2007. The starting date of the sample is chosen because of the relative stability of inflation during this period and to mitigate other issues (for example, a structural break) that may influence our findings.

Three years worth of data are reserved for out-of-sample comparisons of paths one to twelve months into the future. To maintain the forecast window constant, we expand the estimation sample one month at a time until the end of the estimation sample (December 2005).

Figure 3 plots the *MSFE* and the *MSFP* in terms of the percentage forecast improvement of the model that includes money aggregates (M2 in the top panel, M3 in the bottom panel) relative to the model that excludes them. The forecasting model is based on local projections whose lag length is determined automatically by information criteria (Hurvich and Tsai’s (1989) AIC corrected).

Figure 3 is typical of many forecasting exercises, where additional variables are found to improve forecasts in the short-run but not at medium or longer horizons. In fact, both panels show that the *MSFE* quickly deteriorates with the forecast horizon. Meanwhile, the *MSFP* is very stable over the forecast horizon suggesting that monetary aggregates do improve forecasting accuracy over virtually all horizons considered. These gains are nevertheless rather small (at most 3% improvement) and the Hausman test of predictive ability indeed

suggests that the differences are not statistically significant.

7.2 A Macroeconomic Forecasting Exercise

On June 30, 2004, the Federal Open Market Committee (FOMC) raised the federal funds rate (the U.S. key monetary policy rate) from 1% to 1.25% – a level it had not reached since interest rates were last changed from 1.5% to 1.25% on November 6, 2002. For more than a year before the June 30, 2004 change, the Federal Reserve had kept the federal funds rate fixed at 1%. This section examines forecasts of the U.S. economy on the eve of the first in a series of interest rate increases that would culminate two years later, on June 29, 2006, with the federal funds rate at 5.25%.

Our forecast exercise examines U.S. real GDP growth (on a yearly basis, in percentage terms, and seasonally adjusted); inflation (measured by the personal consumption expenditures deflator on a yearly basis, in percentage terms, and seasonally adjusted); the federal funds rate; and the 10 year Treasury Bond rate. All data are measured quarterly (with the federal funds rate and the 10 year T-Bond rate averaged over the quarter) from 1953:II to 2004:II. With these data, we then construct two-year (eight-quarters) ahead forecasts for this system of variables by local projections. The lag length of the projections was automatically selected to be six by AIC_C – a correction to AIC designed for autoregressions and with better properties in small samples than either AIC or SIC (see Hurvich and Tsai, 1989).

Figure 4 displays these forecasts along with the actual realizations of these economic variables, conditional and marginal 95% bands, and 95% Scheffé bands. Several results deserve comment. First, the 95% Scheffé bands are more conservative and tend to fan out as the forecast horizon increases but over the two-year period examined, they tend to be

relatively close to the traditional 95% marginal bands (specially for U.S. GDP). Second, the 95% conditional bands are considerably narrower in all cases but they are meant to capture the uncertainty generated by that period's shock, not the overall uncertainty of the path. Third, our simple experiment results in projections for output and inflation that are more optimistic than the actual data later displayed. As a consequence, our forecast for the federal funds rate is more aggressive (after two years we would have predicted the rate to be at 5.5% instead of 5.25%) although the general pattern of interest rate increases is very similar. Not surprisingly, the 10 year T-Bond rate is also predicted to be higher than it actually was although consistent with a higher inflation premium. For completeness, the same forecasts are displayed in Figure 5 with 95th, 50th and 5th simultaneous percentiles to form appropriate fan charts.

In order to make sense of were the differences between our forecasts and the historical record may come from, we experimented with the following hypothetical. Suppose that the Federal Reserve at the time believed that inflation would not run as high as predicted (perhaps because of the end of major military operations in Iraq suggested more stability in oil markets would be forthcoming or other factors that may be difficult to quantify with the model). Along these lines, we experimented with a path of inflation that tracks the lower 95% conditional confidence band so that inflation is predicted to be at 3.4% (rather than at 3.8%) after two years.

The results of this conditional experiment are reported in Figure 6. We begin by remarking that this hypothetical path is very conservative: the Wald test of the null that the hypothetical is statistically equivalent to the forecast has a p-value of 0.71, that is, the

distance between the hypothetical and the forecasted path is only 29% in probability units. Therefore, we feel reasonably certain that such an experiment is still well approximated by our model. Interestingly, the forecasts obtained by conditioning on this hypothetical path for inflation are remarkably close to the historical record. In particular, the path of increases in the federal funds rate is virtually identical to the actual path whereas the path of the 10 year T-Bond rate is mostly within the 95% conditional bands. The most significant difference was a slight drop in output after one year to a 3% growth rate that in the conditional was predicted to be closer to 3.5%, but otherwise both paths seem to reconnect at the end of the two year predictive horizon. Whereas we cannot be certain that this hypothetical reflected the Federal Reserve view's on inflation at the time, it serves to illustrate that formal statistical experimentation with alternative scenarios can be easily provided to policy makers.

8 Conclusions

Suppose you are comparing two defective computer monitors. One where the color of each pixel is randomly chosen from a relatively tight distribution centered at the true color for each pixel; the other shifts the color of each pixel by a random shock of higher variance than in the first monitor, but by the shock is the same for all pixels. Although the individual pixel, color error-rate in the first monitor is smaller than in the second, images in the first will be quite blurry while images in the second will be perfectly crisp – albeit with the wrong tint.

By the same token, we believe it is more sensible to compare the patterns implied by the forecast paths of competing models jointly when assessing predictive ability. To that end, our paper provides a long list of results. We begin by deriving the asymptotic distribution of

forecast paths generated by finite order VARs or local projections from potentially infinite order DGPs. Hence our results cover a wide class of situations practitioners are likely to encounter in practice, leaving for further research elaborate extensions from these foundational results.

Returning to the intuition of our computer monitor example, we ask what is the best way to compare the predictive ability of models. We accomplish this in two ways: by extending the traditional mean squared forecast error measures to paths (and hence we create the mean squared forecast path) and by extending Diebold-Mariano-West statistics of equal predictive ability in terms on the joint null over the path rather than on its constituent elements.

Summarizing the wealth of information contained in the joint distribution of the forecast path presents a host of new difficulties. We provide several new graphical solutions to this problem based on the observation that the Cholesky decomposition of the covariance matrix of the forecast path orthogonalizes the path into the constituent shocks hitting the forecasting model at each horizon. Hence we introduce simultaneous confidence bands based on Scheffé's (1953) method, conditional bands and fan charts based on the quantiles of the joint predictive density. In the end, the underlying intuition of our derivations is the same as the intuition in classical linear regression with correlated regressors: while individual coefficients may be imprecisely estimated (low t-statistics), the joint effect could still be quite precisely estimated (high F-statistics).

Knowledge of the joint distribution is also advantageous for counterfactual simulation. However, experimentation with alternative scenarios is complicated in most economic applications since the Lucas Critique warns of the possibility that the forecasting model may be

parametrically unstable with respect to the hypothetical path. Furthermore, insofar as the hypothetical paths are far away from the history observed, we are asking an approximate model to make predictions in regions where the model has no training from the data. To get a grip on these issues, we provide formal statistics on the distance of the hypothetical from the average historical distribution of the paths based on the Wald principle. Once the validity of the hypothetical is formally assessed, we provide simple formulas to derive the paths and their distribution, conditional on the hypothetical.

References

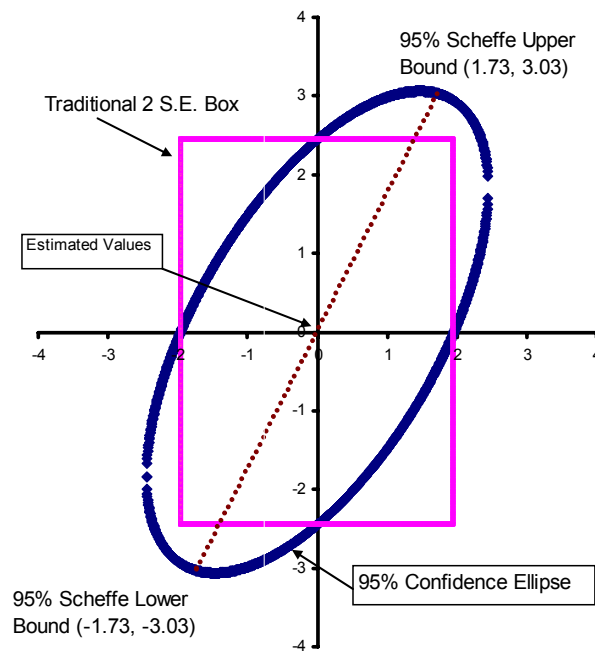
- Anderson, Theodore W. (1994) **The Statistical Analysis of Time Series Data**. New York, New York: Wiley Interscience.
- Bhansali, R. J. (2002) “Multi-Step Forecasting,” in **A Companion to Economic Forecasting**, Michael P. Clements and David F. Hendry (eds.). Oxford, U.K.: Blackwell Publishers.
- Bowden, David C. (1970) “Simultaneous Confidence Bands for Linear Regression Models,” *Journal of the American Statistical Association*, 65(329): 413-421.
- Chao, J.C., Valentina Corradi and Norman R. Swanson (2001) “An Out-of-Sample Test for Granger-Causality,” *Macroeconomic Dynamics*, 5: 598-620.
- Clark, Todd E. and Michael W. McCracken (2001) “Tests of Equal Forecast Accuracy and Encompassing of Nested Models,” *Journal of Econometrics*, 105: 85-110.
- Clements, Michael P. and David F. Hendry (1993) “On the Limitations of Comparing Mean Square Forecast Errors,” *Journal of Forecasting*, 12: 617-676.
- Corradi, Valentina, Norman R. Swanson and Claudio Olivetti (2001) “Predictive Ability with Cointegrated Variables,” *Journal of Econometrics*, 104: 315-358.
- Diebold, Francis X. and Roberto S. Mariano (1995) “Comparing Predictive Accuracy,” *Journal of Business and Economic Statistics*, 13: 253-263.
- Giacomini, Raffaella and Halbert White (2006) “Tests of Conditional Predictive Ability,” *Econometrica*, 74(6): 1545-1578.
- Gonçalves, Silvia and Lutz Kilian (2006) “Asymptotic and Bootstrap Inference for $AR(\infty)$ Processes with Conditional Heteroskedasticity,” *Econometric Reviews*, forthcoming.

- Hamilton, James D. (1994) **Time Series Analysis**. Princeton, NJ: Princeton University Press.
- Horowitz, Joel L. (2001) "The Bootstrap and Hypothesis Tests in Econometrics," *Journal of Econometrics*, 100(1): 37-40.
- Hurvich, Clifford M. and Chih-Ling Tsai (1989) "Regression and Time Series Model Selection in Small Samples," *Biometrika*, 76(2): 297-307.
- Jordà, Òscar (2005) "Estimation and Inference of Impulse Responses by Local Projections," *American Economic Review*, 95(1): 161-182.
- Jordà, Òscar and Sharon Kozicki (2007) "Estimation and Inference by the Method of Projection Minimum Distance," U.C. Davis, Department of Economics, Working Paper 07-8.
- Kuersteiner, Guido (2001) "Optimal Instrumental Variables Estimation for ARMA Models," *Journal of Econometrics*, 104(2): 359-405.
- Kuersteiner, Guido (2002) "Efficient Instrumental Variables Estimation for Autoregressive Models with Conditional Heteroskedasticity," *Econometric Theory*, 18: 547-583.
- Leeper, Eric M. and Tao Zha (2003) "Modest Policy Interventions," *Journal of Monetary Economics*, 50: 1673-1700.
- Lehmann, E. L. and Joseph P. Romano (2005) **Testing Statistical Hypothesis**. Berlin, Germany: Springer-Verlag.
- Lewis, R. A. and Gregory C. Reinsel (1985) "Prediction of Multivariate Time Series by Autoregressive Model Fitting," *Journal of Multivariate Analysis*, 16(33): 393-411.
- Lucas Jr., Robert E. (1976) "Econometric Policy Evaluation: A Critique," in Brunner, Karl and A. H. Meltzer (eds.). *Carnegie-Rochester Conference Series on Public Policy*, 1: 104-130.
- Lütkepohl, Helmut (2005) **New Introduction to Multiple Time Series**. Berlin, Germany: Springer-Verlag.
- Lütkepohl, Helmut and P.S. Poskitt (1991) "Estimating Orthogonal Impulse Responses via Vector Autoregressive Models," *Econometric Theory*, 7: 487-496.
- Marcellino, Massimiliano, James H. Stock and Mark W. Watson (2003) "Macroeconomic Forecasting in the Euro Area: Country-Specific versus Area-Wide Information," *European Economic Review*, 47(1): 1-18.
- Marcellino, Massimiliano, James H. Stock and Mark W. Watson (2006) "A Comparison of Direct and Iterated Multistep AR Methods for Forecasting Macroeconomic Time Series," *Journal of Econometrics*, 127(1-2): 499-526.

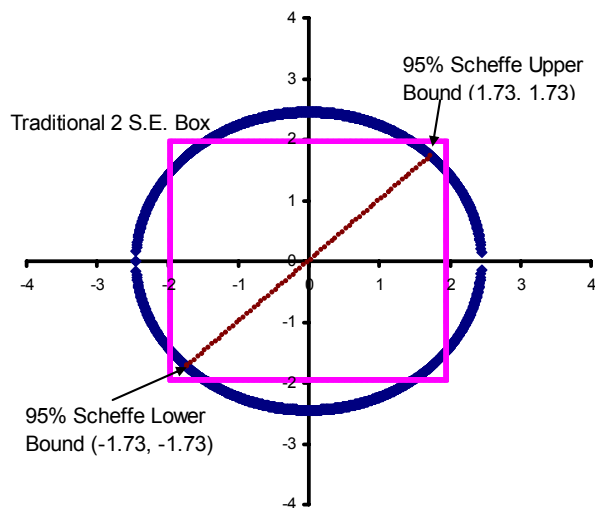
- McCracken, Michael W. (2000) "Robust Out-of-Sample Inference," *Journal of Econometrics*, 99: 195-223.
- Scheffé, Henry (1953) A Method for Judging All Contrasts in the Analysis of Variance," *Biometrika*, 40: 87-104.
- Stock, James H. and Mark W. Watson (1999) "Forecasting Inflation," *Journal of Monetary Economics*, 44(2): 293-335.
- Waggoner, Daniel F. and Tao Zha (1999) "Conditional Forecasts in Dynamic Multivariate Models," *Review of Economics and Statistics*, 81(4): 639-651.
- West, Kenneth D. (1996) "Asymptotic Inference about Predictive Ability," *Econometrica*, 64: 1067-1084.
- White, Halbert (1980) "A Heteroskedasticity-Consistent Covariance Matrix Estimator and a Direct Test of Heteroskedasticity," *Econometrica*, 48(4): 817-838.

Figure 1 – 95% Confidence Ellipse for AR(1) Forecast Path over Two Horizons

Panel 1 – Standard confidence bands, confidence ellipse, and Scheffé Bounds

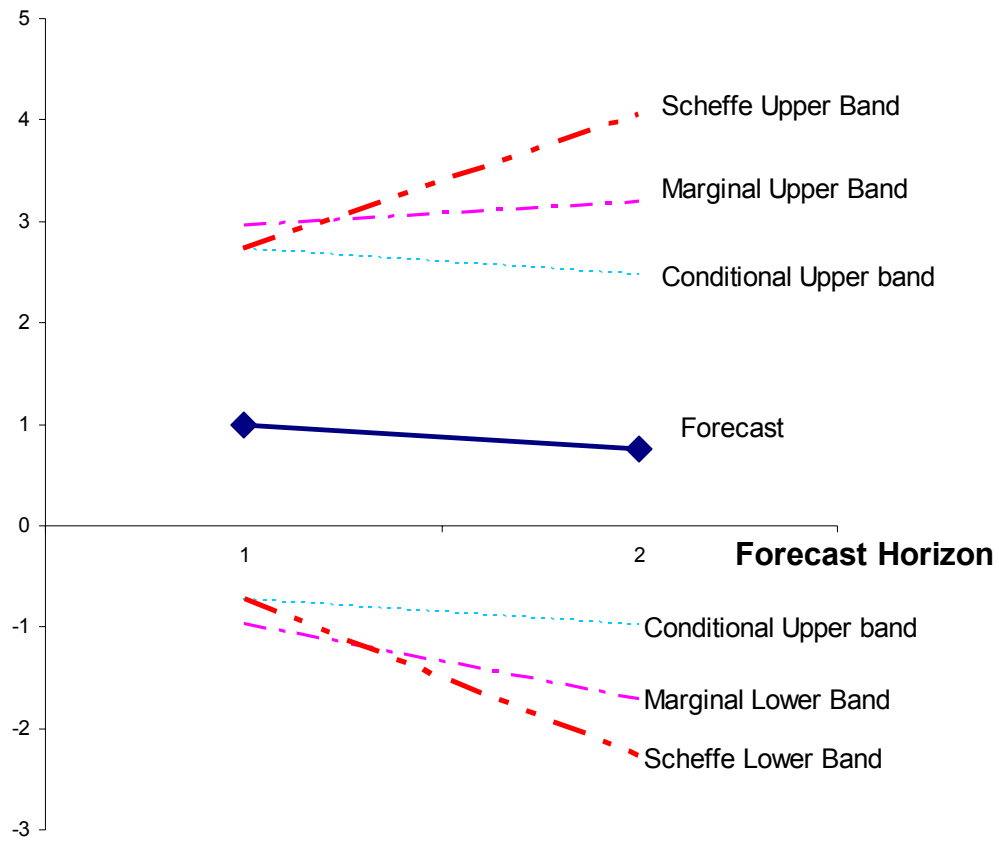


Panel 2 – 95% Confidence Circle for Orthogonalized Forecast Path



Notes: AR Coefficient = 0.75, Error Variance = 1

Figure 2 – 95% Confidence Scheffé, Marginal and Conditional Bands



Notes: AR(1) two period ahead forecasts with $\rho = 0.75$ and $\sigma = 1$. This representation corresponds to the two dimensional representation in figure 1.

Figure 3 – Predicting U.S. Inflation with and without Monetary Aggregates

Sample: January 1985 – January 2007

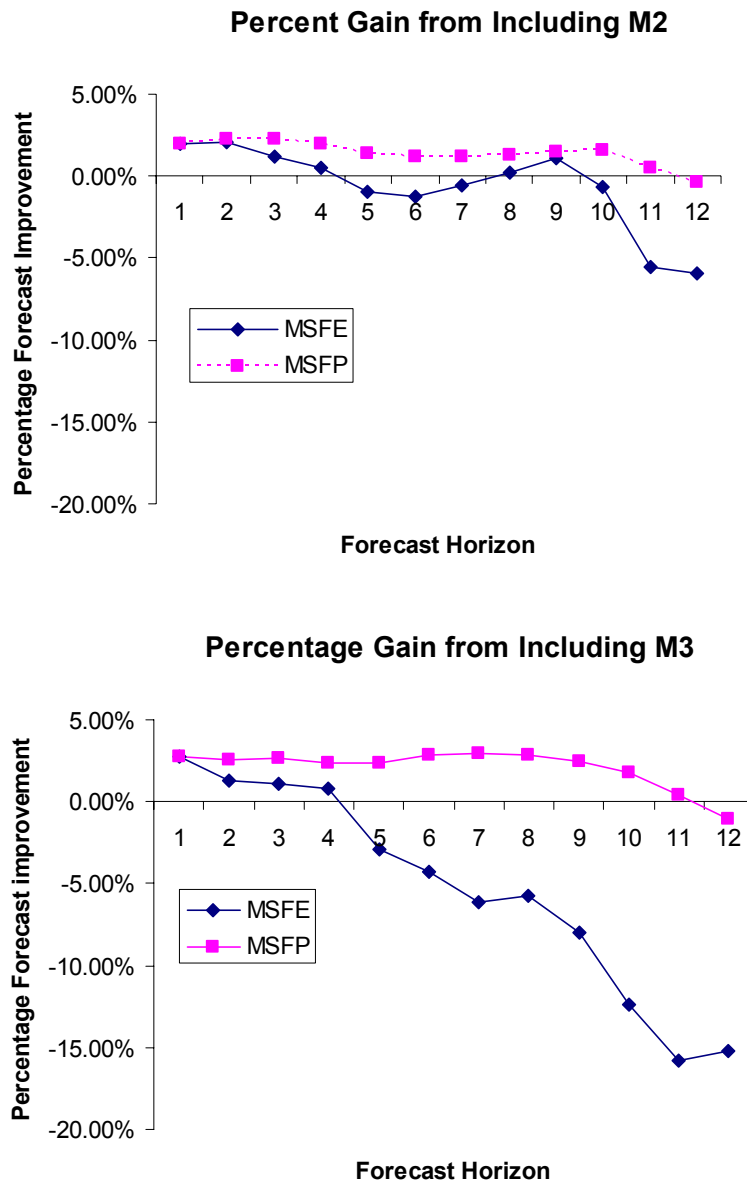
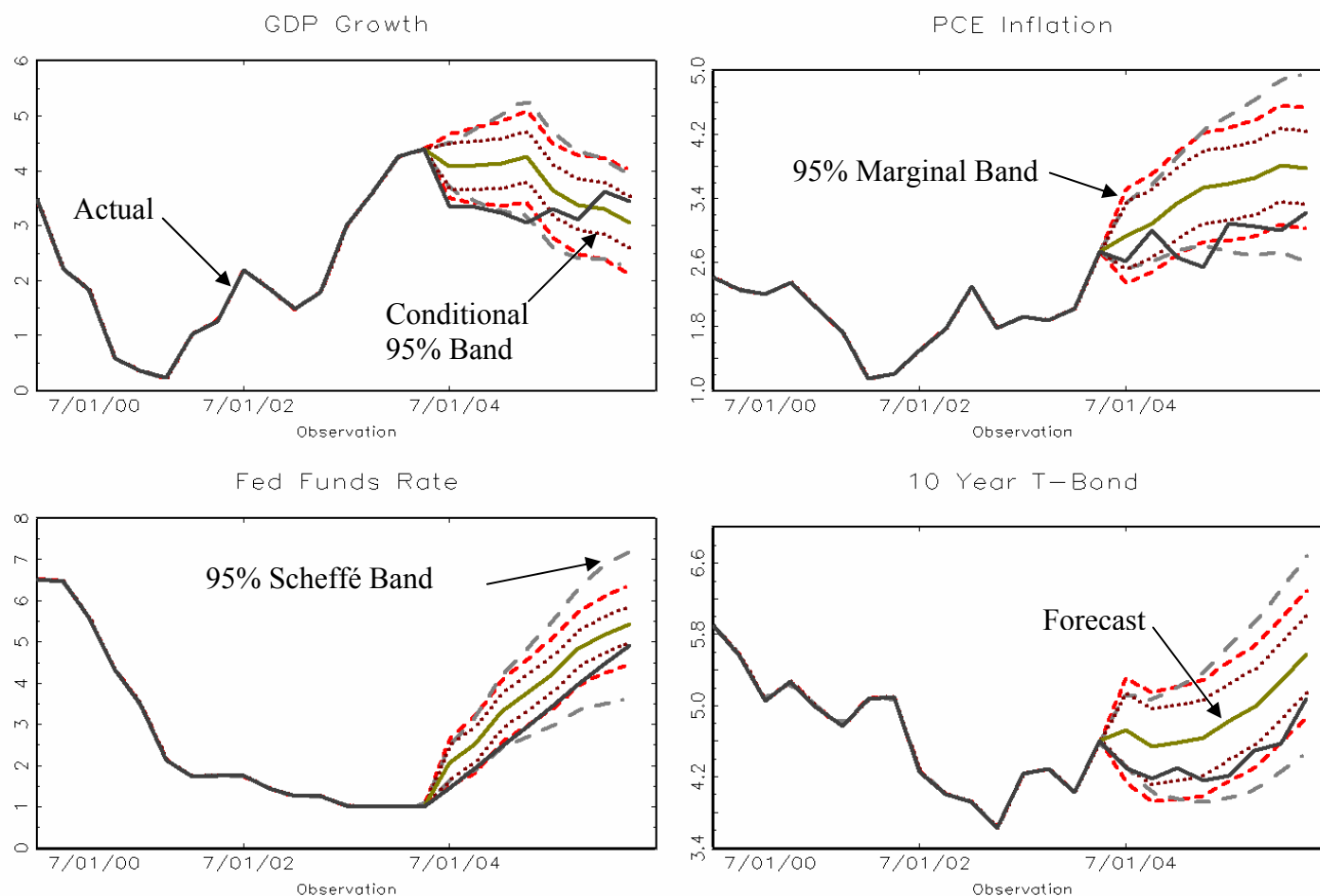
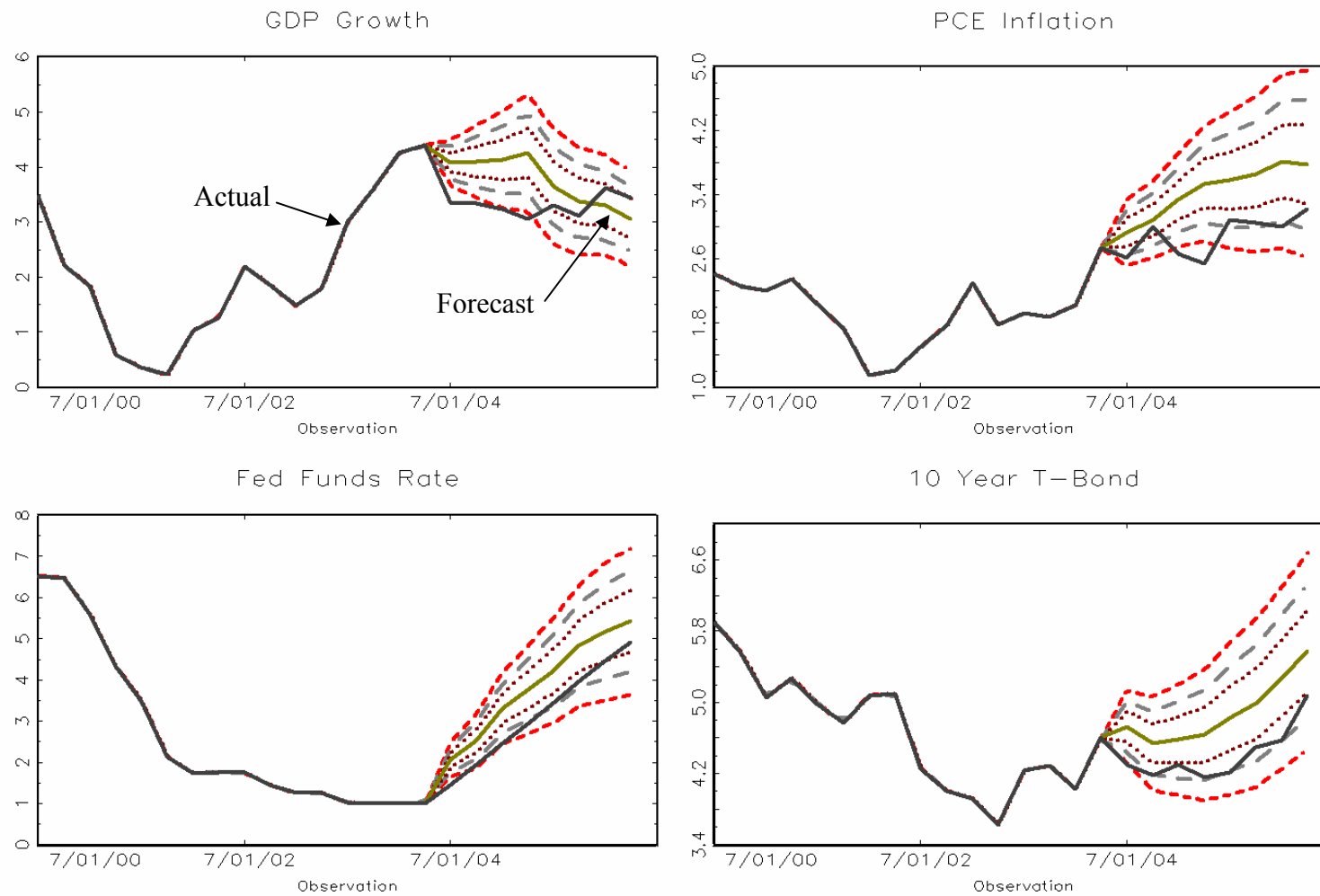


Figure 4 – 95% Scheffé, Conditional and Marginal Error Bands for Macroeconomic Forecasts



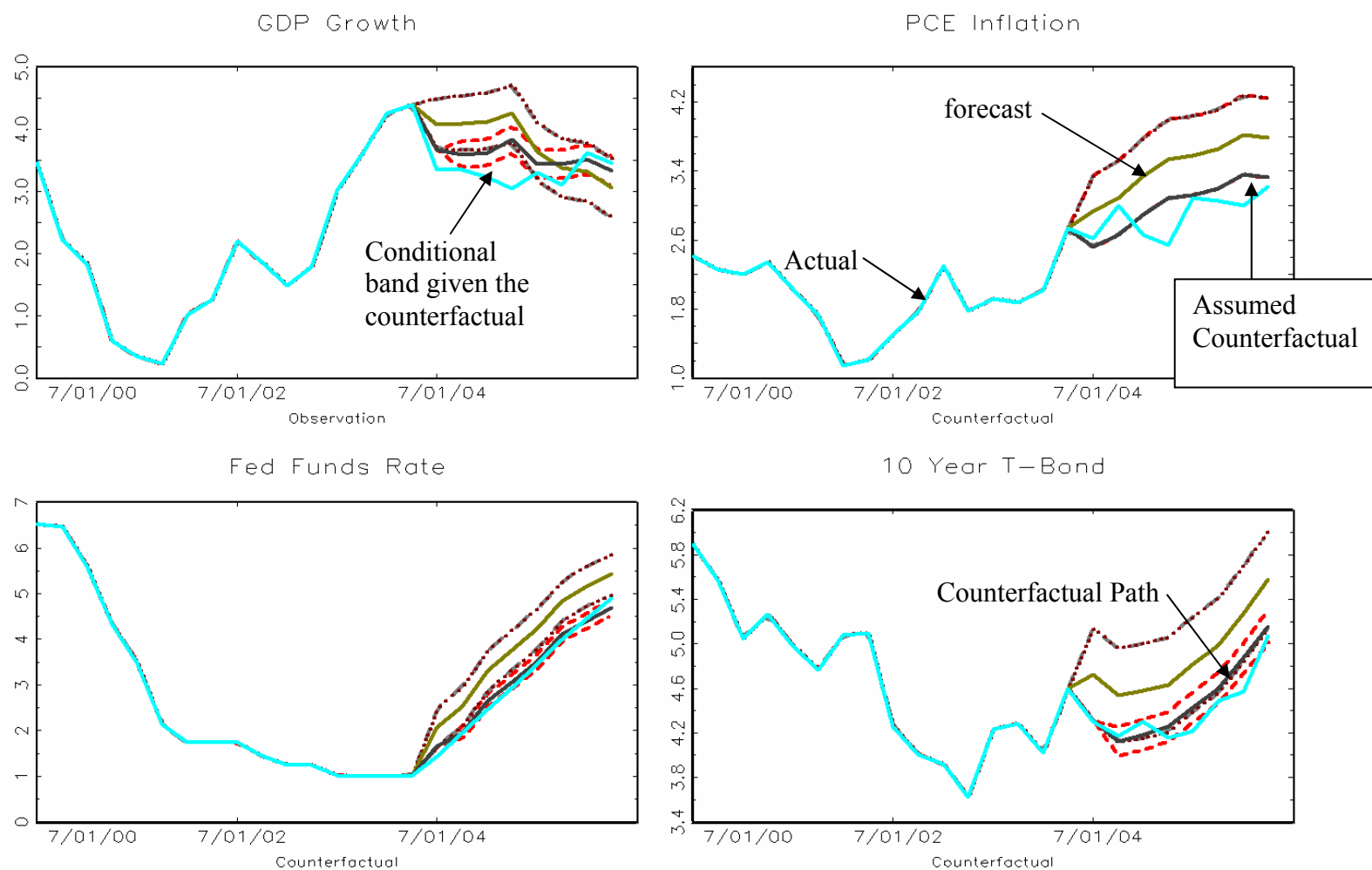
Notes: Sample runs from 1953:II to 2004:II. The forecast horizon runs from 2004:III to 2006:III.

Figure 5 – Fan Charts: 95th, 50th, and 5th percentiles



Notes: These are the same forecasts as in figure 4 but with Scheffé percentiles displayed.

Figure 6 – Counterfactual: Inflation set at the lower 95% conditional band value. Sample: 1953:II – 2004:II



Notes: Distance of the counterfactual from the forecast in probability units is 0.29 (or p-value of the joint test of equality is 0.71).