

FORECAST EVALUATION

Kenneth D. West
University of Wisconsin

January 2005

ABSTRACT

This chapter summarizes recent literature on asymptotic inference about forecasts. Both analytical and simulation based methods are discussed. The emphasis is on techniques applicable when the number of competing models is small. Techniques applicable when a large number of models is compared to a benchmark are also briefly discussed.

Paper prepared for the Handbook of Economic Forecasting. I thank participants in the January 2004 preconference, Pablo M. Pincheira-Brown, Todd E. Clark and Michael W. McCracken for helpful comments. I also thank Pablo M. Pincheira-Brown for research assistance and the National Science Foundation for financial support.

1. INTRODUCTION

This chapter reviews asymptotic methods for inference about moments of functions of predictions and prediction errors. The methods may rely on conventional asymptotics or they may be bootstrap based. The relevant class of applications are ones in which the investigator uses a long time series of predictions and prediction errors as a model evaluation tool. Typically the evaluation is done retrospectively rather than in real time. A classic example is Meese and Rogoff's (1983) evaluation of exchange rate models.

In most applications, the investigator aims to compare two or more models. Measures of relative model quality might include ratios or differences of mean, mean-squared or mean-absolute prediction errors; correlation between one model's prediction and another model's realization (also known as forecast encompassing); or comparisons of utility or profit-based measures of predictive ability. In other applications, the investigator focuses on a single model, in which case measures of model quality might include correlation between prediction and realization, lack of serial correlation in one step ahead prediction errors, ability to predict direction of change, or bias in predictions.

Predictive ability has long played a role in evaluation of econometric models. An early example of a study that retrospectively set aside a large number of observations for predictive evaluation is Wilson (1934, pp307-308). Wilson, who studied monthly price data spanning more than a century, used estimates from the first half of his data to forecast the next twenty years. He then evaluated his model by computing the correlation between prediction and realization.¹ Growth in data and computing power has led to widespread use of similar predictive evaluation techniques, as is indicated by the applications cited below.

To prevent misunderstanding, it may help to stress that the techniques discussed here are probably of little relevance to studies that set aside one or two or a handful of observations for out of sample evaluation. The reader is referred to textbook expositions about confidence intervals around a prediction, or to proposals for simulation methods such as Fair (1980). As well, the paper does not cover

density forecasts. Inference about such forecasts is covered in the handbook chapter by Corradi and Swanson (2004b).

Finally, the paper takes for granted that one wishes to perform out of sample analysis. My purpose is to describe techniques that can be used by researchers who have decided, for reasons not discussed in this chapter, to use a non-trivial portion of their samples for prediction. See recent work by Chen (2004), Clark and McCracken (2005) and Inoue and Kilian (2004a, 2004b) for different takes on the possible power advantages of using out of sample tests.

Section 2 illustrates the evolution of the relevant methodology. For clarity and simplicity, this section, and indeed most of the paper, focuses on tests for equal mean squared prediction error (MSPE). MSPE is not only simple, but it is also arguably the most commonly used measure of predictive ability. Sections 3 through 6 discuss inference when the number of models under evaluation is small. “Small” is not precisely defined, but in sample sizes typically available in economics suggests a number in the single digits. Section 3 discusses inference in the unusual, but conceptually simple, case in which none of the models under consideration rely on estimated regression parameters to make predictions. Sections 4 and 5 relax this assumption, but for reasons described in those sections assume that the models under consideration are nonnested. Section 4 describes when reliance on estimated regression parameters is irrelevant asymptotically, so that section 3 procedures may still be applied. Section 5 describes how to account for reliance on estimated regression parameters. Section 6 considers nested models. Section 7 discusses inference when the number of models being evaluated is large, possibly larger than the sample size. Section 8 concludes.

2. A BRIEF HISTORY

I begin the discussion with a brief history of methodology for inference, focusing on mean squared prediction errors (MSPE).

Let e_{1t} and e_{2t} denote one step ahead prediction errors from two competing models. Let their corresponding second moments be $\sigma_1^2 \equiv Ee_{1t}^2$ and $\sigma_2^2 \equiv Ee_{2t}^2$. (For reasons explained below, the assumption of stationarity—the absence of a t subscript on σ_1^2 and σ_2^2 —is not always innocuous. See below. For the moment, I maintain it for consistency with the literature about to be reviewed.) One wishes to test the null $H_0: \sigma_1^2 - \sigma_2^2 = 0$, or perhaps construct a confidence interval around the point estimate of $\sigma_1^2 - \sigma_2^2$.

Observe that $E(e_{1t} - e_{2t})(e_{1t} + e_{2t}) = \sigma_1^2 - \sigma_2^2$. Thus $\sigma_1^2 - \sigma_2^2 = 0$ if and only if the covariance or correlation between $e_{1t} - e_{2t}$ and $e_{1t} + e_{2t}$ is zero. Let us suppose initially that (e_{1t}, e_{2t}) is i.i.d.. Granger and Newbold (1977) used this observation to suggest testing $H_0: \sigma_1^2 - \sigma_2^2 = 0$ by testing for zero correlation between $e_{1t} - e_{2t}$ and $e_{1t} + e_{2t}$. This procedure was earlier proposed by Morgan (1939) in the context of testing for equality between variances of two normal random variables. Granger and Newbold (1977) assumed that the forecast errors had zero mean, but Morgan (1939) indicates that this assumption is not essential. The Granger and Newbold test was extended to multistep, serially correlated and possibly non-normal prediction errors by Meese and Rogoff (1988) and Mizrahi (1995).

Ashley et al. (1980) proposed a test of equal MSPE in the context of nested models. For nested models, equal MSPE is theoretically equivalent to a test of Granger non-causality. Ashley et al. (1980) proposed executing a standard F-test, but with out of sample prediction errors used to compute restricted and unrestricted error variances. Ashley et al. (1980) recommended that tests be one-sided, testing whether the unrestricted model has smaller MSPE than the restricted (nested) model: it is not clear what it means if the restricted model has a significantly smaller MSPE than the unrestricted model.

The literature on predictive inference that is a focus of this chapter draws on now standard central limit theory introduced into econometrics research by Hansen (1982)—what I will call “standard results” in the rest of the discussion. Perhaps the first explicit use of standard results in predictive inference is Christiano (1989). Let $f_t = e_{1t}^2 - e_{2t}^2$. Christiano observed that we are interested in the mean of f_t , call it $Ef_t \equiv \sigma_1^2 - \sigma_2^2$.² And there are standard results on inference about means—indeed, if f_t is i.i.d. with finite variance,

introductory econometrics texts describe how to conduct inference about Ef_t given a sample of $\{f_t\}$. A random variable like $e_{1t}^2 - e_{2t}^2$ may be non-normal and serially correlated. But results in Hansen (1982) apply to non-i.i.d. time series data. (Details below.)

One of Hansen's (1982) conditions is stationarity. Christiano acknowledged that standard results might not apply to his empirical application because of a possible failure of stationarity. Specifically, Christiano compared predictions of models estimated over samples of increasing size: the first of his 96 predictions relied on models estimated on quarterly data running from 1960 to 1969, the last from 1960 to 1988. Because of increasing precision of estimates of the models, forecast error variances might decline over time. (This is one sense in which the assumption of stationarity was described as "not obviously innocuous" above.)

West et al. (1993) and West and Cho (1995) independently used standard results to compute test statistics. The objects of interest were MSPEs and a certain utility based measure of predictive ability. Diebold and Mariano (1995) proposed using the same standard results, also independently, but in a general context that allows one to be interested in the mean of a general loss or utility function. As detailed below, these papers explained either in context or as a general principle how to allow for multistep, non-normal, and conditionally heteroskedastic prediction errors.

The papers cited in the preceding two paragraphs all proceed without proof. None directly address the possible complications from parameter estimation noted by Christiano (1989). A possible approach to allowing for these complications in special cases is in Hoffman and Pagan (1989) and Ghysels and Hall (1990). These papers showed how standard results from Hansen (1982) can be extended to account for parameter estimation in out of sample tests of instrument residual orthogonality when a fixed parameter estimate is used to construct the test. (Christiano (1989), and most of the forecasting literature, by contrast updates parameter estimate as forecasts progress through the sample.) A general analysis was first presented in West (1996), who showed how standard results can be extended

when a sequence of parameter estimates is used, and for the mean of a general loss or utility function.

Further explication of developments in inference about predictive ability requires me to start writing out some results. I therefore call a halt to the historical summary. The next section (“small number of models, non-nested”) discusses results related to the papers cited here. Section 4 discusses results from a subsequent literature that allows nested models, but still assumes the number of models is small. I then discuss literature that assumes that there is a large number of models.

3. A SMALL NUMBER OF NONNESTED MODELS, PART I

Analytical results are clearest in the unusual (in economics) case in which predictions do not rely on estimated regression parameters. As in the previous section, let e_{1t} and e_{2t} denote prediction errors from two competing models; let $\sigma_1^2 \equiv Ee_{1t}^2$ and $\sigma_2^2 \equiv Ee_{2t}^2$ denote their second moments; and, finally, let $f_t = e_{1t}^2 - e_{2t}^2$. Here and whenever possible, for simplicity and clarity I assume covariance stationarity—neither the first nor second moments of f_t depend on t . At present (predictions do not depend on estimated regression parameters), this assumption is innocuous. It allows simplification of formulas. The results below can be extended to allow moment drift as long as time series averages converge to suitable constants. See Giacomini and White (2003). I do not assume that e_{1t} and e_{2t} are i.i.d. They may be non-normal and serially correlated.

Suppose we have a sample of predictions of size P . Let $\bar{f}^* \equiv P^{-1} \sum_t f_t$ denote the sample mean of f_t , $\bar{f}^* = P^{-1} \sum_t e_{1t}^2 - P^{-1} \sum_t e_{2t}^2$. (The reason for the “*” superscript will become apparent below.) Then the “standard result” referenced above is the time series central limit theorem:

$$(3.1) \quad \sqrt{P}(\bar{f}^* - Ef_t) \sim_A N(0, V^*), \quad V^* \equiv \sum_{j=-\infty}^{\infty} E(f_t - Ef_t)(f_{t-j} - Ef_t).$$

For example, if (e_{1t}, e_{2t}) is i.i.d., then so, too, is $e_{1t}^2 - e_{2t}^2$, and $V^* = E(f_t - Ef_t)^2 = \text{variance}(e_{1t}^2 - e_{2t}^2)$. In such a case, as the number of forecast errors $P \rightarrow \infty$ one can estimate V^* consistently with $\hat{V}^* = P^{-1} \sum_t (f_t - \bar{f}^*)^2$ or (if

the null is imposed in estimation of V^*) $\hat{V}^* = P^{-1} \sum f_t^2$.

Of course, $e_{1t}^2 - e_{2t}^2$ may be serially correlated. There may be serial correlation in second moments (GARCH, for example), even if (e_{1t}, e_{2t}) is a martingale difference sequence. And if (e_{1t}, e_{2t}) is a multistep forecast error, $e_{1t}^2 - e_{2t}^2$ will be serially correlated even if (e_{1t}, e_{2t}) is a moving average of i.i.d. innovations. In such a case, one can estimate V^* with a heteroskedasticity and autocorrelation consistent covariance matrix estimator. In any event, given a consistent estimator \hat{V}^* , one can test the null $H_0: \sigma_1^2 - \sigma_2^2 = 0$ with the t-statistic

$$(3.2) \quad \bar{f}^* / [\hat{V}^*/P]^{1/2} \underset{A}{\sim} N(0,1), \quad \hat{V}^* \underset{p}{\rightarrow} V^* \equiv \sum_{j=-\infty}^{\infty} E(f_t - Ef_t)(f_{t-j} - Ef_t).$$

Confidence intervals can be constructed in obvious fashion from $[\hat{V}^*/P]^{1/2}$. The extension to testing whether all of a set of models have equal MSPE is straightforward, as described in the next section.

To return to the general case: Let \bar{f}^* be the sample moment of interest (possibly a $m \times 1$ vector), \bar{f}^* a sample average of a stationary vector f_t . (In addition to MSPE, moments of interest include, for example, mean error, mean absolute error and mean utility or profit.) Under seemingly weak conditions (see below), \bar{f}^* is asymptotically normal with variance-covariance matrix

$$(3.3) \quad V^* \equiv \sum_{j=-\infty}^{\infty} E(f_t - Ef_t)(f_{t-j} - Ef_t)'$$

Let \hat{V}^* be a consistent estimator of V^* . If the null is $Ef_t = 0$, a Wald test, computed as

$$(3.4) \quad \bar{f}^{*'} \hat{V}^{*-1} \bar{f}^*,$$

is asymptotically $\chi^2(m)$ under the null.

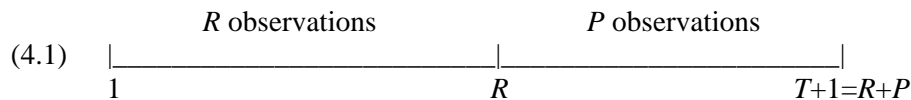
It is well known that asymptotically normality of \bar{f}^* is consistent with quite general forms of dependence and heterogeneity in the forecast errors (White (1984)). The “seemingly” in the reference to “seemingly weak” assumptions in the previous paragraph refers to some ancillary assumptions that are

not satisfied in some relevant applications. First, the number of models m must be “small” relative to the number of predictions P . In an extreme case in which $m > P$, conventional estimators will yield \hat{V}^* that is not of full rank. As well, and more informally, one suspects that conventional asymptotics will yield a poor approximation if m is large relative to P . Section 7 briefly discusses alternative approaches likely to be useful in such contexts.

Second, and more generally, V^* must be full rank. When the number of models $m=2$, this rules out $e_{1t}^2 = e_{2t}^2$ with probability 1 (obviously). It also rules out pairs of models in which $\sqrt{P}(\hat{\sigma}_1^2 - \hat{\sigma}_2^2) \rightarrow_p 0$. This latter condition is violated in applications in which one or both models make predictions based on estimated regression parameters, and the models are nested. This is discussed in section 6 below.

4. A SMALL NUMBER OF NONNESTED MODELS, PART II

In the vast majority of economic applications, one or more of the models under consideration rely on estimated regression parameters when making predictions. To spell out the implications for inference, it is necessary to define some additional notation. For simplicity, assume that one step ahead prediction errors are the object of interest. Let the total sample size be $T+1$. The last P observations of this sample are used for forecast evaluation. The first R observations are used to construct an initial set of regression estimates that are then used for the first prediction. We have $R+P=T+1$. Schematically:



In the forecasting literature, three distinct schemes figure prominently in how one generates the sequence of regression estimates necessary to make predictions. Asymptotic results differ slightly for the three, so it is necessary to distinguish between them. Let β denote the vector of regression parameters whose estimates are used to make predictions. In the *recursive* scheme, the size of the sample used to

estimate β grows as one makes predictions for successive observations. One first estimates β with data from 1 to R and uses the estimate to predict observation $R+1$ (recall that I am assuming one step ahead predictions, for simplicity); one then estimates β with data from 1 to $R+1$, with the new estimate used to predict observation $R+2$;; finally, one estimate β with data from 1 to T , with the final estimate used to predict observation $T+1$. In the *rolling* scheme, the sequence of β 's is always generated from a sample of size R . The first estimate of β is obtained with a sample running from 1 to R , the next with a sample running from 2 to $R+1$, ..., the final with a sample running from $T-R+1$ to T . In the *fixed* scheme, one estimates β just once, using data from 1 to R . In all three schemes, the number of predictions is P and the size of the smallest regression sample is R . Examples of applications using each of these schemes include Faust et al. (2004) (recursive), Cheung et al. (2003) (rolling) and Ashley et al. (1980) (fixed). The fixed scheme is relatively attractive when it is computationally difficult to update parameter estimates. The rolling scheme is relatively attractive when one wishes to guard against moment or parameter drift that is difficult to model explicitly.

It may help to illustrate with a simple example. Suppose one model under consideration is a univariate zero mean AR(1): $y_t = \beta^* y_{t-1} + e_{1t}$. Then the sequence of P estimates of β^* are generated as follows for $t=R, \dots, T$:

$$(4.2) \quad \begin{aligned} \text{recursive: } \hat{\beta}_t &= [\sum_{s=1}^t (y_{s-1}^2)]^{-1} [\sum_{s=1}^t y_{s-1} y_s]; \\ \text{rolling: } \hat{\beta}_t &= [\sum_{s=t-R+1}^t (y_{s-1}^2)]^{-1} [\sum_{s=t-R+1}^t y_{s-1} y_s]; \\ \text{fixed: } \hat{\beta}_t &= [\sum_{s=1}^R (y_{s-1}^2)]^{-1} [\sum_{s=1}^R y_{s-1} y_s]. \end{aligned}$$

In each case, the one step ahead prediction error is $\hat{e}_{t+1} \equiv y_{t+1} - y_t \hat{\beta}_t$. Observe that for the fixed scheme $\hat{\beta}_t = \hat{\beta}_R$ for all t , while $\hat{\beta}_t$ changes with t for the rolling and recursive schemes.

Suppose there are two models, say $y_t = X_{1t}' \beta_1^* + e_{1t}$ and $y_t = X_{2t}' \beta_2^* + e_{2t}$. (Two parenthetical comments: (1) Note the dating convention: X_{1t} and X_{2t} can be used to predict y_t , for example $X_{1t} = y_{t-1}$ if model 1 is an

AR(1). (2)The assumption of linearity is made for expositional convenience. As noted below, the relevant asymptotic theory allows nonlinear models and estimators, including GMM and maximum likelihood.) The population MSPEs are $\sigma_1^2 \equiv Ee_{1t}^2$ and $\sigma_2^2 \equiv Ee_{2t}^2$. (Absence of a subscript t on the MSPEs is for simplicity and without substance.) Define the sample one step ahead forecast errors and sample MSPEs as

$$(4.4) \quad \hat{e}_{1t+1} \equiv y_{t+1} - X_{1t+1}' \hat{\beta}_{1t}, \hat{e}_{2t+1} \equiv y_{t+1} - X_{2t+1}' \hat{\beta}_{2t}, \hat{\sigma}_1^2 = P^{-1} \sum_{t=R}^T \hat{e}_{1t+1}^2, \hat{\sigma}_2^2 = P^{-1} \sum_{t=R}^T \hat{e}_{2t+1}^2.$$

With MSPE the object of interest, one examines the difference between the sample MSPEs $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$. Let

$$(4.5) \quad \hat{f}_t \equiv \hat{e}_{1t}^2 - \hat{e}_{2t}^2, \bar{f} \equiv P^{-1} \sum_{t=R}^T \hat{f}_{t+1} \equiv \hat{\sigma}_1^2 - \hat{\sigma}_2^2.$$

Observe that \bar{f} defined in (4.5) differs from \bar{f}^* defined above (3.1) in that \bar{f} relies on \hat{e} 's, whereas \bar{f}^* relies on e 's.

The null hypothesis is $\sigma_1^2 - \sigma_2^2 = 0$. One way to test the null would be to substitute \hat{e}_{1t} and \hat{e}_{2t} for e_{1t} and e_{2t} in the formulas presented in the previous section. If $(e_{1t}, e_{2t})'$ is i.i.d., for example, one would set $\hat{V}^* = P^{-1} \sum_{t=R}^T (\hat{f}_{t+1} - \bar{f})^2$, compute the t-statistic

$$(4.6) \quad \bar{f} / [\hat{V}^* / P]^{1/2}$$

and use standard normal critical values. (I use the “*” in \hat{V}^* for both $P^{-1} \sum_{t=R}^T (\hat{f}_{t+1} - \bar{f})^2$ [this section] and for $P^{-1} \sum_{t=R}^T (f_{t+1} - \bar{f})^2$ [previous section] because under the asymptotic approximations described below, both are consistent for the long run variance of f_{t+1} .)

Use of (4.6) is not obviously an advisable approach. Clearly, $\hat{e}_{1t}^2 - \hat{e}_{2t}^2$ is polluted by error in estimation of β_1 and β_2 . It is not obvious that sample averages of $\hat{e}_{1t}^2 - \hat{e}_{2t}^2$ (i.e., \bar{f}) have the same asymptotic distribution as those of $e_{1t}^2 - e_{2t}^2$ (i.e., \bar{f}^*). Under suitable conditions (see below), a key factor determining whether the asymptotic distributions are equivalent is whether or not the two models are nested. If they

are nested, the distributions are not equivalent. Use of (4.6) with normal critical values is not advised. This is discussed in the next section.

If the models are not nested, West (1996) showed that when conducting inference about MSPE, parameter estimation error is *asymptotically irrelevant*. I put the phrase in italics because I will have frequent recourse to it in the sequel: “asymptotic irrelevance” means that one conduct inference by applying standard results to the mean of the loss function of interest, treating parameter estimation error as irrelevant.

To explain this result, as well as to illustrate when asymptotic irrelevance does not apply, requires some—actually, considerable—notation. I will phase in some of this notation in this section, with most of the algebra deferred to the next section. Let β^* denote the $k \times 1$ population value of the parameter vector used to make predictions. Suppose for expositional simplicity that the model(s) used to make predictions are linear,

$$y_t = X_t' \beta^* + e_t$$

if there is a single model,

$$y_t = X_{1t}' \beta_1^* + e_{1t}, y_t = X_{2t}' \beta_2^* + e_{2t}, \beta^* \equiv (\beta_1^{*'}, \beta_2^{*'})'$$

if there are two competing models. Let $f_t(\beta^*)$ be the random variable whose expectation is of interest.

Then leading examples of $f_t(\beta^*)$ include:

$$(4.7a) \quad f_t(\beta^*) = e_{1t}^2 - e_{2t}^2 = (y_t - X_{1t}' \beta_1^*)^2 - (y_t - X_{2t}' \beta_2^*)^2, (Ef_t = 0 \text{ means equal MSPE});$$

$$(4.7b) \quad f_t(\beta^*) = e_t = y_t - X_t' \beta^* (Ef_t = 0 \text{ means zero mean prediction error});$$

$$(4.7c) \quad f_t(\beta^*) = e_{1t} X_{2t}' \beta_2^* = (y_t - X_{1t}' \beta_1^*) X_{2t}' \beta_2^* (Ef_t = 0 \text{ means zero correlation between one model's prediction error and another model's prediction, an implication of forecast encompassing proposed by Chong and Hendry (1986)});$$

$$(4.7d) \quad f_t(\beta^*) = e_{1t}(e_{1t} - e_{2t}) = (y_t - X_{1t}' \beta_1^*) [(y_t - X_{1t}' \beta_1^*) - (y_t - X_{2t}' \beta_2^*)] (Ef_t = 0 \text{ is an implication of forecast encompassing proposed by Harvey et al. (1998)});$$

$$(4.7e) \quad f_t(\beta^*) = e_{t+1}e_t = (y_{t+1}-X_{t+1}'\beta^*)(y_t-X_t'\beta^*) \quad (Ef_t=0 \text{ means zero first order serial correlation});$$

$$(4.7f) \quad f_t(\beta^*) = e_tX_t'\beta^* = (y_t-X_t'\beta^*)X_t'\beta^* \quad (Ef_t=0 \text{ means the prediction and prediction error are uncorrelated});$$

$$(4.7g) \quad f_t(\beta^*) = |e_{1t}| - |e_{2t}| = |y_t-X_{1t}'\beta_1^*| - |y_t-X_{2t}'\beta_2^*|, \quad (Ef_t=0 \text{ means equal mean absolute error}).$$

More generally, $f_t(\beta^*)$ can be per period utility or profit, or differences across models of per period utility or profit, as in Leitch and Tanner (1991) or West et al. (1993).

Let $\hat{f}_{t+1} \equiv f_{t+1}(\hat{\beta}_t)$ denote the sample counterpart of $f_{t+1}(\beta^*)$, with $\bar{f} \equiv P^{-1}\sum_{t=R}^T \hat{f}_{t+1}$ the sample mean evaluated at the series of estimates of β^* . Let $\bar{f}^* \equiv P^{-1}\sum_{t=R}^T f_{t+1}(\beta^*)$ denote the sample mean evaluated at β^* . Let F denote the $(1 \times k)$ derivative of the expectation of f_t , evaluated at β^* :

$$(4.8) \quad F = \partial Ef_t(\beta^*)/\partial \beta.$$

For example, $F = -EX_t'$ for mean prediction error (4.7b).

Then under mild conditions,

$$(4.9) \quad \sqrt{P(\bar{f}-Ef_t)} = \sqrt{P(\bar{f}^*-Ef_t)} + F \times (P/R)^{1/2} \times [O_p(1) \text{ terms from the sequence of estimates of } \beta^*] + o_p(1).$$

Some specific formulas are in the next section. Result (4.9) holds not only when f_t is a scalar, as I have been assuming, but as well when f_t is a vector. (When f_t is a vector of dimension (say) m , F has dimension $m \times k$.)

Thus, uncertainty about the estimate of Ef_t can be decomposed into uncertainty that would be present even if β^* were known and, possibly, additional uncertainty due to estimation of β^* . The qualifier “possibly” results from at least three sets of circumstances in which error in estimation of β^* is asymptotically irrelevant: (1) $F=0$; (2) $P/R \rightarrow 0$; (3) the variance of the terms due to estimation of β^* is exactly offset by the covariance between these terms and $\sqrt{P(\bar{f}^*-Ef_t)}$. For cases (1) and (2), the middle term in (4.9) is identically zero ($F=0$) or vanishes asymptotically ($P/R \rightarrow 0$), implying that $\sqrt{P(\bar{f}-Ef_t)} - \sqrt{P(\bar{f}^*-Ef_t)} \rightarrow_p 0$; for case (3) the asymptotic variances of $\sqrt{P(\bar{f}-Ef_t)}$ and $\sqrt{P(\bar{f}^*-Ef_t)}$ happen to be the same. In any of the

three sets of circumstances, inference can proceed as described in the previous section. This is important because it simplifies matters if one can abstract from uncertainty about β^* when conducting inference.

To illustrate each of the three circumstances:

1. For MSPE in our linear example $F = (-2EX_{1t}'e_{1t}, 2EX_{2t}'e_{2t})'$. So $F=0_{1 \times k}$ if the predictors are uncorrelated with the prediction error.³ Similarly, $F=0$ for mean absolute prediction error (4.7g) ($E[|e_{1t}| - |e_{2t}|]$) when the prediction errors have a median of zero, conditional on the predictors. (To prevent confusion, let me emphasize that MSPE and mean absolute error are unusual in that asymptotic irrelevance applies even when P/R is not small. In this sense, my focus on MSPE is a bit misleading.)

Let me illustrate the implications with an example in which f_t is a vector rather than a scalar. Suppose that we wish to test equality of MSPEs from $m+1$ competing models, under the assumption that the forecast error vector $(e_{1t}, \dots, e_{m+1,t})'$ is i.i.d.. Define the $m \times 1$ vectors

$$(4.10) \quad f_t \equiv (e_{1t}^2 - e_{2t}^2, \dots, e_{1t}^2 - e_{m+1,t}^2)', \quad \hat{f}_t = (\hat{e}_{1t}^2 - \hat{e}_{2t}^2, \dots, \hat{e}_{1t}^2 - \hat{e}_{m+1,t}^2)', \quad \bar{f} = P^{-1} \sum_{t=R}^T \hat{f}_{t+1}.$$

The null is that $E f_t = 0_{m \times 1}$. (Of course, it is arbitrary that the null is defined as discrepancies from model 1's squared prediction errors; test statistics are identical regardless of the model used to define f_t .) Then under the null

$$(4.11) \quad \bar{f}' \hat{V}^* \bar{f} \sim_A \chi^2(m), \quad \hat{V}^* \rightarrow_p V^* \equiv \sum_{j=-\infty}^{\infty} E(f_t - E f_t)(f_{t-j} - E f_t)',$$

where, as indicated, \hat{V}^* is a consistent estimate of the $m \times m$ long run variance of f_t . If $f_t \equiv (e_{1t}^2 - e_{2t}^2, \dots, e_{1t}^2 - e_{m+1,t}^2)'$ is serially uncorrelated (sufficient for which is that $(e_{1t}, \dots, e_{m+1,t})'$ is i.i.d.), then a possible estimator of V is simply

$$\hat{V}^* = P^{-1} \sum_{t=R}^T (\hat{f}_{t+1} - \bar{f})(\hat{f}_{t+1} - \bar{f})'.$$

If the squared forecast errors display persistence (GARCH and all that), a robust estimator of the variance-covariance matrix should be used (West and Cho (1995)).

2. One can see in (4.9) that asymptotic irrelevance holds quite generally when $P/R \rightarrow 0$. The intuition is that the relatively large sample (big R) used to estimate β produces small uncertainty relative to uncertainty that would be present in the relatively small sample (small P) even if one knew β . The result was noted informally by Chong and Hendry (1986). Simulation evidence in West (1996, 2001), McCracken (2004) and Clark and McCracken (2001) suggests that $P/R < .10$ more or less justifies using the asymptotic approximation that assumes asymptotic irrelevance.

3. This fortunate cancellation of variance and covariance terms occurs for certain moments and loss functions, when estimates of parameters needed to make predictions are generated by the recursive scheme (but not by the rolling or fixed schemes), and when forecast errors are conditionally homoskedastic. These loss functions are: mean prediction error; serial correlation of one step ahead prediction errors; zero correlation between one model's forecast error and another model's forecast. This is illustrated in point 5b in section 6 below.

To repeat: When asymptotic irrelevance applies, one can proceed as described in the previous section. One need not account for dependence of forecasts on estimated parameter vectors. When asymptotic irrelevance does not apply, matters are more complicated. This is discussed in the next two sections.

5. A SMALL NUMBER OF NONNESTED MODELS, PART III

Asymptotic irrelevance fails in a number of important cases, at least according to the asymptotics of West (1996). Under the rolling and fixed schemes, it fails quite generally. For example, it fails for mean prediction error, correlation between realization and prediction, encompassing, and zero correlation in one step ahead prediction errors. Under the recursive scheme, it similarly fails for such moments when prediction errors are not conditionally homoskedastic. In such cases, asymptotic inference requires accounting for uncertainty about parameters used to make predictions.

To motivate the decomposition (4.9), I will temporarily switch from my example of comparison of MSPEs to one in which one is looking at mean prediction error. The variable f_t is thus redefined to equal the prediction error, $f_t = e_t$, and Ef_t is the moment of interest. I will further use a trivial example, in which the only predictor is the constant term, $y_t = \beta^* + e_t$. Let us assume as well, as in the Hoffman and Pagan (1989) and Ghysels and Hall (1990) analysis of predictive tests of instrument-residual orthogonality, that the fixed scheme is used and predictions are made using a single estimate of β^* . This single estimate is the least squares estimate on the sample running from 1 to R , $\hat{\beta}_R = R^{-1} \sum_{s=1}^R y_s$. Now, $\hat{e}_{t+1} = e_{t+1} - (\hat{\beta}_R - \beta^*) = e_{t+1} - R^{-1} \sum_{s=1}^R e_s$. So

$$(5.1) \quad P^{-1/2} \sum_{t=R}^T \hat{e}_{t+1} = P^{-1/2} \sum_{t=R}^T e_{t+1} - (P/R)^{1/2} (R^{-1/2} \sum_{s=1}^R e_s).$$

This is in the form (4.9), with: $F = -1$, $R^{-1/2} \sum_{s=1}^R e_s = [O_p(1)$ terms due to the sequence of estimates of β^*], and the $o_p(1)$ term identically zero.

The asymptotic approximations discussed in this section assume a large sample of both predictions and prediction errors,

$$(5.2) \quad P \rightarrow \infty, R \rightarrow \infty, \lim_{T \rightarrow \infty} \frac{P}{R} = \pi, 0 \leq \pi < \infty.$$

If e_t is well behaved, say i.i.d. with finite variance σ^2 , the bivariate vector $(P^{-1/2} \sum_{t=R}^T e_{t+1}, R^{-1/2} \sum_{s=1}^R e_s)'$ is asymptotically normal with variance covariance matrix $\sigma^2 I_2$. It follows that

$$(5.3) \quad P^{-1/2} \sum_{t=R}^T e_{t+1} - (P/R)^{1/2} (R^{-1/2} \sum_{s=1}^R e_s) \sim_A N(0, (1+\pi)\sigma^2).$$

Thus use of $\hat{\beta}_R$ rather than β^* in predictions inflates the asymptotic variance of the estimator of mean prediction error by a factor of $1+\pi$. Clearly the sample analogue to π is P/R . If P/R is small, the implied value of π is small and inflation factor is small; otherwise the inflation factor can be big.

To write out the result more generally: Let h_t be the orthogonality condition used to identify β^* ,

$$(5.4a) \quad h_t = X_t e_t$$

or

$$(5.4b) \quad h_t = (X_{1t}' e_{1t}, X_{2t}' e_{2t})'$$

in our least squares example (where $h_t = h_t(\beta^*)$, with β^* suppressed for simplicity). (More generally, h_t is the score if the estimation method is maximum likelihood, or the GMM orthogonality condition.) Let the average of h_t that figures into estimation of $\hat{\beta}_t$ be called $H(t)$. In our linear least squares example,

$$(5.5) \quad H(t) = t^{-1} \sum_{s=1}^t h_s \text{ (recursive), } H(t) = R^{-1} \sum_{s=t-R+1}^t h_s \text{ (rolling), } H(t) = R^{-1} \sum_{s=1}^R h_s \text{ (fixed),}$$

$$\bar{H} = P^{-1} \sum_{t=R}^T H(t) \text{ (recursive, rolling and fixed).}$$

Suppose that the estimator of β^* can be written $\hat{\beta}_t - \beta^* = B(t)H(t) + \text{terms of small magnitude}$, and let B be the large sample counterpart of $B(t)$. B is the inverse of the expectation of the Hessian (ML) or linear combination of orthogonality conditions (GMM). In our linear least squares example, B is $k \times k$,

$$(5.6a) \quad B = (EX_1 X_1')^{-1}$$

or the $k \times k$ block diagonal matrix

$$(5.6b) \quad B = \text{diag}[(EX_1 X_1')^{-1}, (EX_2 X_2')^{-1}].$$

For generality, allow f_t to possibly be a vector, with f_t having dimension $m \times 1$, $m \geq 1$. Write the long run variance of $(f_{t+1}', h_t)'$ as

$$(5.7) \quad S = \begin{bmatrix} V^* & S_{fh} \\ S_{fh}' & S_{hh} \end{bmatrix}$$

Here, $V^* \equiv \sum_{j=-\infty}^{\infty} E(f_t - E f_t)(f_{t-j} - E f_{t-j})'$ is $m \times m$, $S_{fh} \equiv \sum_{j=-\infty}^{\infty} E(f_t - E f_t)h_{t-j}'$ is $m \times k$, and $S_{hh} \equiv \sum_{j=-\infty}^{\infty} E h_t h_{t-j}'$ is $k \times k$, and f_t and h_t are understood to be evaluated at β^* . The asymptotic ($R \rightarrow \infty$) variance-covariance matrix of the estimator of β^* is

$$(5.8) \quad V_{\beta} \equiv B S_{hh} B'.$$

With π defined in (5.2), define the scalars λ_{fh} , λ_{hh} and $\lambda \equiv (1 + \lambda_{hh} - 2\lambda_{fh})$

(5.9)	Sampling scheme	λ_{fh}	λ_{hh}	λ
	recursive	$1 - \pi^{-1} \ln(1 + \pi)$	$2[1 - \pi^{-1} \ln(1 + \pi)]$	1
	rolling, $\pi \leq 1$	$\frac{\pi}{2}$	$\pi \frac{\pi^2}{3}$	$1 - \frac{\pi^2}{3}$
	rolling, $\pi > 1$	$1 - \frac{1}{2\pi}$	$1 - \frac{1}{3\pi}$	$\frac{2}{3\pi}$
	fixed	0	π	$1 + \pi$

Finally, define the $m \times k$ matrix F as in (4.8), $F \equiv \partial E f_t(\beta^*) / \partial \beta$.

Then under conditions such as those described below, an expansion yields

$$(5.10) \quad P^{-1/2} [\sum_{t=R}^T f(\hat{\beta}_{t+1}) - E f_t] = P^{-1/2} [\sum_{t=R}^T f_{t+1}(\beta^*) - E f_t] + F(P/R)^{1/2} [BR^{1/2} \bar{H}] + o_p(1)$$

(Equation (5.10) is written in a form to make connection with (5.1) apparent. The connection with (4.9)

may be more apparent if we rewrite (5.10) as $\sqrt{P}(\bar{f} - E f_t) = \sqrt{P}(\bar{f}^* - E f_t) + F(P/R)^{1/2} [BR^{1/2} \bar{H}] + o_p(1)$. As well,

$P^{-1/2} [\sum_{t=R}^T f(\hat{\beta}_{t+1}) - E f_t]$ is asymptotically normal with variance-covariance matrix

$$(5.11) \quad V = V^* + \lambda_{fh} (FBS_{fh}' + S_{fh} B' F') + \lambda_{hh} FV_{\beta} F'$$

V^* is the long run variance of $P^{-1/2} [\sum_{t=R}^T f_{t+1}(\beta^*) - E f_t]$ and is the same object as V^* defined in (3.3),

$\lambda_{hh} FV_{\beta} F'$ is the long run variance of $F(P/R)^{1/2} [BR^{1/2} \bar{H}]$, and $\lambda_{fh} (FBS_{fh}' + S_{fh} B' F')$ is the covariance between the two.

When uncertainty about β^* matters asymptotically, the adjustment to the standard error that would be appropriate if predictions were based on population rather than estimated parameters is increasing in:

(1) the ratio of number of predictions P to number of observations in smallest regression sample R (note

that as $\pi \rightarrow 0$, $\lambda_{fh} \rightarrow 0$ and $\lambda_{hh} \rightarrow 0$) and (2) the variance-covariance matrix of the estimator of the parameters

used to make predictions. Both conditions are intuitive. Simulations in West (1996, 2001), West and

McCracken (1998), McCracken (2000), Chao et al. (2001) and Clark and McCracken (2001, 2003)

indicate that with plausible parameterizations for P/R and uncertainty about β^* , failure to adjust the

standard error can result in very substantial size distortions. It is possible that $V < V^*$ – that is, accounting for uncertainty about regression parameters may *lower* the asymptotic variance of the estimator.⁴ This happens in some leading cases of practical interest, when the rolling scheme is used. See point 6b in section 6 for an illustration.

A consistent estimator of V results from using the obvious sample analogues. A possibility is to compute λ_{fh} and λ_{hh} from (5.9) setting $\pi=P/R$. (See Table 1 for the implied formulas for λ_{fh} , λ_{hh} and λ .) As well, one can estimate F from the sample average of $\partial f(\hat{\beta}_t)/\partial\beta$, $\hat{F} = P^{-1}\sum_{t=R}^T \partial f(\hat{\beta}_t)/\partial\beta$;⁵ estimate V_β and B from one of the sequence of estimates of β^* . For example, for mean prediction error, for the fixed scheme, one might set

$$\hat{F} = P^{-1}\sum_{t=R}^T X_{t+1}', \hat{B} = (R^{-1}\sum_{s=1}^R X_s X_s')^{-1}, \hat{V}_\beta \equiv (R^{-1}\sum_{s=1}^R X_s X_s')^{-1}(R^{-1}\sum_{s=1}^R X_s X_s' \hat{e}_s^2)(R^{-1}\sum_{s=1}^R X_s X_s')^{-1}.$$

Here, \hat{e}_s , $1 \leq s \leq R$ is the in-sample least squares residual associated with the parameter vector $\hat{\beta}_R$ that is used to make predictions and the formula for \hat{V}_β is the usual heteroskedasticity consistent covariance matrix for $\hat{\beta}_R$. (Other estimators are also consistent, for example sample averages running from 1 to T .) Finally, one can combine these with an estimate of the long run variance S constructed using a heteroskedasticity and autocorrelation consistent covariance matrix estimator (Newey and West (1987, 1994), Andrews (1991), Andrews and Monahan (1994), den Haan and Levin (2000)).

Alternatively, one can compute a smaller dimension long run variance as follows. Let us assume for the moment that f_t and hence V are scalar. Define the (2×1) vector \hat{g}_t as

$$(5.12) \quad \hat{g}_t = \begin{bmatrix} \hat{f}_t \\ \hat{F}\hat{B}\hat{h}_t \end{bmatrix}.$$

Let g_t be the population counterpart of \hat{g}_t , $g_t \equiv (f_t, FBh_t)'$. Let Ω be the (2×2) long run variance of g_t , $\Omega \equiv \sum_{j=-\infty}^{\infty} E g_t g_{t-j}'$. Let $\hat{\Omega}$ be an estimate of Ω . Let $\hat{\Omega}_{ij}$ be the (i,j) element of $\hat{\Omega}$. Then one can consistently estimate V with

$$(5.13) \quad \hat{V} = \hat{\Omega}_{11} + 2\lambda_{fh}\hat{\Omega}_{12} + \lambda_{hh}\hat{\Omega}_{22}.$$

The generalization to vector f_t is straightforward. Suppose f_t is say $m \times 1$ for $m \geq 1$. Then

$$g_t = \begin{bmatrix} f_t \\ FBh_t \end{bmatrix}$$

is $2m \times 1$, as is \hat{g}_t ; Ω and $\hat{\Omega}$ are $2m \times 2m$. One divides $\hat{\Omega}$ into four ($m \times m$) blocks, and computes

$$(5.14) \quad \hat{V} = \hat{\Omega}(1,1) + \lambda_{fh}[\hat{\Omega}(1,2) + \hat{\Omega}(2,1)] + \lambda_{hh}\hat{\Omega}(2,2).$$

In (5.14), $\hat{\Omega}(1,1)$ is the $m \times m$ block in the upper left hand corner of $\hat{\Omega}$, $\hat{\Omega}(1,2)$ is the $m \times m$ block in the upper right hand corner of $\hat{\Omega}$, and so on.

Alternatively, in some common problems, and if the models are linear, regression based tests can be used. By judicious choice of additional regressors (as suggested for in-sample tests by Pagan and Hall (1983), Davidson and McKinnon (1984) and Wooldridge (1990)), one can “trick” standard regression packages into computing standard errors that properly reflect uncertainty about β^* . See West and McCracken (1998) and Table 3 below for details, Hueng and Wong (2000), Avramov (2002) and Ferreira (2004) for applications.

Conditions for the expansion (5.10) and the central limit result (5.11) include the following.

- Parametric models and estimators of β are required. Similar results may hold with nonparametric estimators, but, if so, these have yet to be established. Linearity is not required. One might be basing predictions on nonlinear time series models, for example, or restricted reduced forms of simultaneous equations models estimated by GMM.
- At present, results with $I(1)$ data are restricted to linear models (Corradi et al. (2001), Rossi (2003)). The results of the previous section continue to apply when $F=0$ or $\pi=0$. When those conditions fail, however, the normalized estimator of Ef_t typically is no longer asymptotically normal. (By $I(1)$ data, I mean $I(1)$ data entered in levels in the regression model. Of course, if one induces stationarity by taking differences or imposing cointegrating relationships prior to estimating β^* , the theory in the present section is

applicable quite generally.)

- Condition (5.2) holds. The next section discusses implications of an alternative asymptotic approximation due to Giacomini and White (2003) that holds R fixed.
- For the recursive scheme, condition (5.2) can be relaxed to allow $\pi=\infty$, with the same asymptotic approximation. (Recall that π is the limiting value of P/R .) Since $\pi<\infty$ is required for rolling and fixed, researchers using those schemes should treat the asymptotic approximation with extra caution if $P\gg R$.
- The expectation of the loss function f must be differentiable in a neighborhood of β^* . This rules out direction of change as a loss function.
- A full rank condition on the long run variance of $(f_{t+1}', (Bh_t)')'$. A necessary condition is that the long run variance of f_{t+1} is full rank. For MSPE, and i.i.d. forecast errors, this means that the variance of $e_{1t}^2 - e_{2t}^2$ be positive (note the absence of a “^” over e_{1t}^2 and e_{2t}^2). This condition will fail in applications in which the models are nested, for in that case $e_{1t} \equiv e_{2t}$, as discussed in the next section. Of course, for the sample forecast errors, $\hat{e}_{1t} \neq \hat{e}_{2t}$ (note the “^”) because of sampling error in estimation of β_1^* and β_2^* . So the failure of the rank condition may not be apparent in practice. McCracken’s (2004) analysis of nested models shows that under the conditions of the present section apart from the rank condition, $\sqrt{P}(\hat{\sigma}_1^2 - \hat{\sigma}_2^2) \rightarrow_p 0$. The next section discusses inference for predictions from such nested models.

6. SMALL NUMBER OF MODELS, NESTED

Analysis of nested models per se does not invalidate the results of the previous section. A rule of thumb is: if the rank of the data becomes degenerate when regression parameters are set at their population values, then a rank condition assumed in the previous section likely is violated. When only two models are being compared, “degenerate” means identically zero.

Consider, as an example, out of sample tests of Granger causality (e.g., Stock and Watson (1999, 2002)). In this case model 2 might be a bivariate VAR, model 1 a univariate AR that is nested in model 2

by imposing suitable zeroes in the model 2 regression vector. If the lag length is 1, for example:

$$(6.1a) \text{ Model 1: } y_t = \beta_{10} + \beta_{11}y_{t-1} + e_{1t} \equiv X_{1t}'\beta_1^* + e_{1t}, X_{1t} \equiv (1, y_{t-1})', \beta_1^* \equiv (\beta_{10}, \beta_{11})';$$

$$(6.1b) \text{ Model 2: } y_t = \beta_{20} + \beta_{21}y_{t-1} + \beta_{22}x_{t-1} + e_{2t} \equiv X_{2t}'\beta_2^* + e_{2t}, X_{2t} \equiv (1, y_{t-1}, x_{t-1})', \beta_2^* \equiv (\beta_{20}, \beta_{21}, \beta_{22})'.$$

Under the null of no Granger causality from x to y , $\beta_{22}=0$ in model 2. Model 1 is then nested in model 2.

Under the null, then,

$$\beta_2^* = (\beta_1^*, 0), X_{1t}'\beta_1^* = X_{2t}'\beta_2^*,$$

and the disturbances of model 2 and model 1 are identical: $e_{2t}^2 - e_{1t}^2 = 0$, $e_{1t}(e_{1t} - e_{2t}) = 0$ and $|e_{1t}| - |e_{2t}| = 0$ for all t .

So the theory of the previous section does not apply if MSPE, $\text{cov}(e_{1t}, e_{1t} - e_{2t})$ or mean absolute error is the moment of interest. On the other hand, the random variable $e_{1t}x_{t-1}$ is nondegenerate under the null, so one can use the theory of the previous section to examine whether $Ee_{1t}x_{t-1} = 0$. Indeed, Chao et al. (2001) show that (5.10) and (5.11) apply when testing $Ee_{1t}x_{t-1} = 0$ with out of sample prediction errors.

The remainder of this section considers the implications of a test that does fail the rank condition of the theory of the previous section—specifically, MSPE in nested models. This is a common occurrence in papers on forecasting asset prices, which often use MSPE to test a random walk null against models that use past data to try to predict changes in asset prices. It is also a common occurrence in macro applications, which, as in example (6.1), compare univariate to multivariate forecasts. In such applications, the asymptotic results described in the previous section will no longer apply. In particular, and under the technical conditions of that section (apart from the rank condition), $\sqrt{P}(\hat{\sigma}_1^2 - \hat{\sigma}_2^2) \rightarrow_p 0$. See Clark and McCracken (2001) and McCracken (2004), who derive the asymptotic distribution of $P(\hat{\sigma}_1^2 - \hat{\sigma}_2^2)$ and certain related quantities. (Note that the normalizing factor is the prediction sample size P rather than the usual \sqrt{P} .) They write test statistics as functionals of Brownian motion. For multistep predictions, limiting distributions are not asymptotically free of nuisance parameters. For one step ahead predictions and conditionally homoskedastic data, these papers supply critical values that can be used to perform

hypothesis tests about equal MSPE (see below).

Analytical and simulation evidence in Clark and McCracken (2001, 2003), Clark and West (2004) and McCracken (2004) indicates that in MSPE comparisons in nested models the usual statistic (4.6) is highly non-normal. Use of standard critical values usually results in *far* too few rejections. As well, the usual statistic has very poor power. For both size and power, the usual statistic performs worse the larger the number of irrelevant regressors included in model 2.

To illustrate the sources of these results, consider the following simple example. The two models are:

(6.2) Model 1: $y_t = e_t$; Model 2: $y_t = \beta^* x_t + e_t$; $\beta^* = 0$; e_t a martingale difference sequence with respect to past y 's and x 's.

In (6.2), all variables are scalars. I use x_t instead of X_{2t} to keep notation relatively uncluttered. For concreteness, one can assume $x_t = y_{t-1}$, but that is not required. I write the disturbance to model 2 as e_t rather than e_{2t} because the null (equal MSPE) implies $\beta^* = 0$ and hence that the disturbance to model 2 is identically equal e_t . Nonetheless, for clarity and emphasis I use the "2" subscript for the sample forecast error from model 2, $\hat{e}_{2t+1} = y_{t+1} - x_t \hat{\beta}_t$. In a finite sample, the model 2 sample forecast error differs from the model 1 forecast error, which is simply y_{t+1} . The model 1 and model 2 MSPEs are

$$(6.3) \quad \hat{\sigma}_1^2 \equiv P^{-1} \sum_{t=R}^T y_{t+1}^2, \quad \hat{\sigma}_2^2 \equiv P^{-1} \sum_{t=R}^T \hat{e}_{2t+1}^2 \equiv P^{-1} \sum_{t=R}^T (y_{t+1} - x_t \hat{\beta}_t)^2$$

Since

$$\hat{f}_{t+1} \equiv y_{t+1}^2 - (y_{t+1} - x_t \hat{\beta}_t)^2 = 2y_{t+1} x_t \hat{\beta}_t - (x_t \hat{\beta}_t)^2$$

we have

$$(6.4) \quad \bar{f} \equiv \hat{\sigma}_1^2 - \hat{\sigma}_2^2 = 2(P^{-1} \sum_{t=R}^T y_{t+1} x_t \hat{\beta}_t) - [P^{-1} \sum_{t=R}^T (x_t \hat{\beta}_t)^2].$$

Now,

$$- [P^{-1} \sum_{t=R}^T (x_t \hat{\beta}_t)^2] \leq 0$$

and under the null ($y_{t+1} = e_{t+1} \sim \text{i.i.d.}$)

$$2(P^{-1} \sum_{t=R}^T y_{t+1} x_t \hat{\beta}_t) \approx 0.$$

So under the null we expect

$$(6.5) \quad \bar{f} \equiv \hat{\sigma}_1^2 - \hat{\sigma}_2^2 < 0$$

or: we expect the sample MSPE from the null model to be *less* than that from the alternative model.

The intuition will be unsurprising to those familiar with forecasting. If the null is true, the alternative model introduces noise into the forecasting process: the alternative model attempts to estimate parameters that are zero in population. In finite samples, use of the noisy estimate of the parameter will *raise* the estimated MSPE of the alternative model relative to the null model. So if the null is true, the model 1 MSPE should be smaller by the amount of estimation noise.

The consequences for inferences are explored in Clark and McCracken (2001, 2003), McCracken (2004) and Clark and West (2004). Let me use the simulation results in latter paper for illustration. Following Ashley et al. (1980), one tailed tests were used. That is, the null of equal MSPE is rejected at (say) the 10 percent level only if the alternative model predicts better than model 1:

$$(6.6) \quad \bar{f} / [\hat{V}^* / P]^{1/2} = (\hat{\sigma}_1^2 - \hat{\sigma}_2^2) / [\hat{V}^* / P]^{1/2} > 1.65,$$

$$\hat{V}^* = \text{estimate of long run variance of } \hat{\sigma}_1^2 - \hat{\sigma}_2^2,$$

$$\text{say, } \hat{V}^* = P^{-1} \sum_{t=R}^T (\hat{f}_{t+1} - \bar{f})^2 = P^{-1} \sum_{t=R}^T [\hat{f}_{t+1} - (\hat{\sigma}_1^2 - \hat{\sigma}_2^2)]^2 \text{ if } e_t \text{ is i.i.d..}$$

Since (6.6) is motivated by an asymptotic approximation in which $\hat{\sigma}_1^2 - \hat{\sigma}_2^2$ is centered around zero, we see from (6.5) that the test will tend to be undersized (reject too infrequently). Across 48 sets of simulations, with DGPs calibrated to match key characteristics of asset price data, Clark and West (2004) found that

the median size of a nominal 10% test using the standard result (6.6) was less than 1%. The size was better with bigger R and worse with bigger P . (Some alternative procedures (described below) had median sizes of 8%-13%.) The power of tests using “standard results” was poor: rejection of about 9%, versus 50%-80% for alternatives.⁶

Similar non-normality results apply if one compares forecasting ability using an encompassing test suggested by Harvey et al. (1998). In this test, which is (4.7d) above, $f_t = e_{1t}(e_{1t} - e_{2t})$, with model 1 the restricted model that is nested in model 2. The relevant sample moment is $\bar{f} = P^{-1} \sum_{t=R}^T \hat{e}_{1t+1} (\hat{e}_{1t+1} - \hat{e}_{2t+1})$. Non-normality also applies if one normalizes differences in MSPEs by the unrestricted MSPE to produce an out of sample F-test. See Clark and McCracken (2001, 2003), and McCracken (2004) for analytical and simulation evidence of marked departures from normality. Of various variants of these tests, Clark and McCracken (2004) find that power is best using the Harvey et al. (1998) version of an encompassing test, normalized by unrestricted variance. So for those who use a non-normal test, Clark and McCracken (2001) recommend the statistic that they call “Enc-new:”

$$(6.7) \quad \text{Enc-new} = \bar{f} = \frac{P^{-1} \sum_{t=R}^T \hat{e}_{1t+1} (\hat{e}_{1t+1} - \hat{e}_{2t+1})}{\hat{\sigma}_2^2}, \quad \hat{\sigma}_2^2 \equiv P^{-1} \sum_{t=R}^T \hat{e}_{2t+1}^2.$$

For one step ahead predictions in conditionally homoskedastic models, Clark and McCracken (2001) provide asymptotic critical values for (6.7); McCracken (2004) does the same for MSPE. These values vary with the number of additional regressors in the unrestricted model (in model 2). For multistep forecasts or conditionally homoskedastic data, the theory in Clark and McCracken (2001) and McCracken (2004) does not yield statistics that are asymptotically free of nuisance parameters. Simulations or bootstrapping are necessary.

Viable approaches to testing in nested models include the following (with the first two summarizing the previous paragraph):

(1) Use critical values from Clark and McCracken (2001) and McCracken (2004), for one step ahead

predictions with conditionally homoskedastic disturbances (e.g., Lettau and Ludvigson (2001)).

(2) Simulate/bootstrap your own standard errors (e.g., Mark (1995), Sarno et al. (2004)).

(3) If $P/R \rightarrow 0$, Clark and McCracken (2001) and McCracken (2004) show that asymptotic irrelevance applies. So for small P/R , use standard critical values (e.g., Clements and Galvao (2003)). Simulations in various papers suggest that it generally does little harm to ignore effects from estimation of regression parameters if $P/R \leq 0.1$. Of course, this cutoff is arbitrary. For some data, a larger value is appropriate, for others a smaller value.

(4) For MSPE and one step ahead forecasts, use the standard test if it rejects: if the standard test rejects, a properly sized test most likely will as well (e.g., Shintani (2004)).⁷

(5) If the null model is a trivial one (i.e., $y_t = e_t$), and if the rolling or fixed regression scheme is used, one can account for regression parameters estimated in model 2 as follows: adjust the difference in MSPEs in a certain way, and use standard results as described in section 4 above. (Clark and West (2004)).

(6) Swear off MSPE, substituting a procedure that is asymptotically normal.

(6a) Test whether model 1's prediction error is uncorrelated with model 2's predictors or the subset of model 2's predictors not included in model 1 (Chao et al. (2001)), $f_t = e_{1t} X_{2t}'$ in our linear example or $f_t = e_{1t} x_{t-1}$ in example (6.1). When both models use estimated parameters for prediction (in contrast to (6.2), in which model 1 does not rely on estimated parameters), the Chao et al. (2001) procedure requires adjusting the variance-covariance matrix for parameter estimation error, as described in the previous section.

(6b) If $\beta_2^* \neq 0$, apply an encompassing test in the form (4.7c), $0 = E e_{1t} X_{2t}' \beta_2^*$, (but *not* (4.7d)). To interpret the condition $\beta_2^* \neq 0$, recall that in nested models, the encompassing null means (in our linear example) $X_{2t}' \beta_2^* = X_{1t}' \beta_1^*$. Thus $\beta_2^* = 0$ if and only if $\beta_1^* = 0$, i.e., if the null model is a driftless martingale difference. Such a null is used in (6.2) above and is common in exchange rate applications. Otherwise $\beta_2^* \neq 0$. One will have $\beta_2^* \neq 0$ even in standard tests of a random walk model of asset prices if the null model includes

deterministic terms to reflect nonzero expected returns.

With computation and technical conditions similar to those in West and McCracken (1998), it may be shown that when $\bar{f} = P^{-1} \sum_{t=R}^T \hat{e}_{1t+1} X_{2t+1}' \hat{\beta}_{2t}$, $\beta_2^* \neq 0$, and the models are nested, then

$$(6.8) \quad \sqrt{P} \bar{f} \sim_A N(0, V), \quad V \equiv \lambda V^*, \quad \lambda \text{ defined in (5.9), } V^* \equiv \sum_{j=-\infty}^{\infty} E e_t e_{t-j} (X_{2t}' \beta_2^*) (X_{2t-j}' \beta_2^*).$$

If e_t is an i.i.d. one step ahead forecast error, then $V^* = E e_t^2 (X_{2t}' \beta_2^*)^2$; one can in this case consistently estimate V^* via $P^{-1} \sum_{t=R}^T (\hat{e}_{1t+1} X_{2t+1}' \hat{\beta}_{2t})^2$. For multistep forecasts, one can estimate V^* using standard methods to compute a long run variance. In any event, once an estimate of V^* is obtained, one multiplies the estimate by λ to obtain an estimate of the asymptotic variance of $\sqrt{P} \bar{f}$. Alternatively, one divides the t-statistic by $\sqrt{\lambda}$.⁸

Observe that $\lambda=1$ for the recursive scheme: this is an example in which there is the cancellation of variance and covariance terms noted in point 3 at the end of section 4. For the fixed scheme, $\lambda > 1$, with λ increasing in P/R . So uncertainty about parameter estimates inflates the variance, with the inflation factor increasing in the ratio of the size of the prediction to regression sample. Finally, for the rolling scheme $\lambda < 1$. So use of (6.8) will result in *smaller* standard errors and larger t-statistics than would use of a statistic that ignores the effect of uncertainty about β^* . The magnitude of the adjustment to standard errors and t-statistics is increasing in the ratio of the size of the prediction to regression sample.

(6c) If $\beta_2^* = 0$, and if the rolling or fixed (but *not* the recursive) scheme is used, apply the encompassing test just discussed, setting $\bar{f} = P^{-1} \sum_{t=R}^T e_{1t+1} X_{2t+1}' \hat{\beta}_{2t}$. Note that in contrast to the discussion just completed, there is no “^” over e_{1t+1} : because model 1 is nested in model 2, $\beta_2^* = 0$ means $\beta_1^* = 0$, so $e_{1t+1} = y_{t+1}$ and e_{1t+1} is observable. One can use standard results—asymptotic irrelevance applies. The factor of λ that appears in (6.8) resulted from estimation of β_1^* , and is now absent. So $V = V^*$; if, for example, e_{1t} is i.i.d., one can consistently estimate V with $\hat{V} = P^{-1} \sum_{t=R}^T (e_{1t+1} X_{2t+1}' \hat{\beta}_{2t})^2$.⁹

(6d) If the rolling or fixed regression scheme is used, construct a conditional rather than unconditional test

(Giacomini and White (2003)). This paper makes both methodological and substantive contributions. The methodological contributions are twofold. First, the paper explicitly allows data heterogeneity (e.g., slow drift in moments). This seems to be a characteristic of much economic data. Second, while the paper's conditions are broadly similar to those of the work cited above, its asymptotic approximation holds R fixed while letting $P \rightarrow \infty$.

The substantive contribution is also twofold. First, the objects of interest are moments of \hat{e}_{1t} and \hat{e}_{2t} rather than e_t . (Even in nested models, \hat{e}_{1t} and \hat{e}_{2t} are distinct because of sampling error in estimation of regression parameters used to make forecasts.) Second, and related, the moments of interest are conditional ones, say $E(\hat{\sigma}_1^2 - \hat{\sigma}_2^2 | \text{lagged } y\text{'s and } x\text{'s})$. The Giacomini and White (2003) framework allows general conditional loss functions, and may be used in nonnested as well as nested frameworks.

Let me close with a summary. An expansion and application of the recommendations of the preceding four sections is given in Tables 2 and 3. The rows of Table 2 are organized by sources of critical values. The first row is for tests that rely on standard results. As described in sections 3 and 4, this means that asymptotic normal critical values are used without explicitly taking into account uncertainty about regression parameters used to make forecasts. The second row is for tests that rely on asymptotic normality, but only after adjusting for such uncertainty as described in section 5 and in some of the final points of this section. The third row is for tests for which it would be ill-advised to use asymptotic normal critical values, as described in this section.

The panels of Table 3 are organized by class of application, panel A for a single model, panel B for a pair of nonnested models, panel C for a pair of nested models. Within each panel, rows are organized by the moment being studied.

Tables 2 and 3 aim to make specific recommendations. While the tables are self-explanatory, some qualifications should be noted. First, the rule of thumb that asymptotic irrelevance applies when $P/R < 0.1$ (point A1 in Table 2, note to Table 3A) is just a rule of thumb. Second, as noted in section 4,

asymptotic irrelevance for MSPE or mean absolute error (point A2 in Table 2, B1 and B2 in Table 3) requires that the prediction error is uncorrelated with the predictors (MSPE) or that the disturbance is symmetric conditional on the predictors (mean absolute error). Otherwise, one will need to account for uncertainty about parameters used to make predictions. Third, some of the results in A3 and A4 (Table 2) and the regression results in Table 3A, rows 1-3, and Table 3B, row 3, have yet to be noted. They are established in West and McCracken (1998). Finally, the suggestion to run a regression on a constant and compute a HAC t-stat (e.g., Table 3, panel B, row 1) is just one way to operationalize a recommendation to use standard results. This recommendation is given in non-regression form in equations (4.5) and (4.6) above.

7. LARGE NUMBER OF MODELS

Sometimes an investigator will wish to compare a large number of models. There is no precise definition of large. But for sample of size typical in economics research, procedures in this section probably have limited appeal when the number of models is say in the single digits, and have a great deal of appeal when the number of models is into double digits or above. White's (2000) empirical example examined 3654 models using a sample of size 1560.

I divide the discussion into (A) applications in which there is a natural null model, and (B) applications in which there is no natural null.

(A) Sometimes one has a natural null, or benchmark, model, which is to be compared to an array of competitors. The leading example is a martingale difference model for an asset price, to be compared to a long list of methods claimed in the past to help predict returns. An obvious problem is controlling size.

White's (2000) "reality check" is a bootstrap method for construction of p-values. It assumes asymptotic irrelevance ($P \ll R$) [though the actual asymptotic condition requires $P/R \rightarrow 0$ at a sufficiently

rapid rate (White (2000, p1105))). The basic mechanics are as follows:

(1)Generate prediction errors, using the scheme of choice (recursive, rolling, fixed).

(2)Generate a series of bootstrap samples as follows.

(a)Sample with replacement from the prediction errors. There is no need to generate bootstrap samples of parameters used for prediction because asymptotic irrelevance is assumed to hold. The bootstrap generally needs to account for possible dependency of the data. White (2000) recommends the stationary bootstrap of Politis and Romano (1994).

(b)In each bootstrap sample, compute the object of interest (say, MSPE).

(3)Compute a p-value by comparing the best of the alternative models to the distribution in the bootstrap data.

While White (2000) motivates the method for its ability to tractably handle situations where the number of models is large relative to sample size, the method can be used in applications with a small number of models as well (e.g., Hong and Lee (2003)).

White's (2000) results have stimulated the development of similar procedures. Corradi and Swanson (2004a) indicate how to account for parameter estimation error, when asymptotic irrelevance does not apply. Corradi, Swanson and Olivetti (2001) present extensions to cointegrated environments. Romano and Wolf (2003) propose that test statistics be studentized, to better exploit the benefits of bootstrapping. Hansen (2003) suggests an alternative formulation that has better power when testing for superior, rather than equal, predictive ability.

(B)Sometimes there is no natural null. McCracken and Sapp (2004) propose that one gauge the "false discovery rate" of Storey (2002). That is, one should control the fraction of rejections that are due to type I error. Hansen et al. (2004) propose constructing a set of models that contain the best forecasting model with prespecified asymptotic probability.

8. CONCLUSIONS

This paper has summarized some recent work about inference about forecasts. The emphasis has been on the effects of uncertainty about regression parameters used to make forecasts, when one is comparing a small number of models. Results applicable for a comparison of a large number of models were also discussed. One of the highest priorities for future work is development of asymptotically normal or otherwise nuisance parameter free tests for equal MSPE or mean absolute error in a pair of nested models. At present only special case results are available.

FOOTNOTES

1. Which, incidentally and regrettably, turned out to be negative.
2. Actually, Christiano looked at root mean squared prediction errors, testing whether $\sigma_1 - \sigma_2 = 0$. For clarity and consistency with the rest of my discussion, I cast his analysis in terms of MSPE.
3. Of course, one would be unlikely to forecast with a model that *a priori* is expected to violate this condition, though prediction is sometimes done with realized right hand side endogenous variables (e.g., Meese and Rogoff (1983)). But prediction exercises do sometimes find that this condition does not hold. That is, out of sample prediction errors display correlation with the predictors (even though in sample residuals often display zero correlation by construction). So even for MSPE, one might want to account for parameter estimation error when conducting inference.
4. Mechanically, such a fall in asymptotic variance indicates that the variance of terms resulting from estimation of β^* is more than offset by a negative covariance between such terms and terms that would be present even if β^* were known.
5. See McCracken (2000) for an illustration of estimation of F for a non-differentiable function.
6. Note that (4.6) and the left hand side of (6.6) are identical, but that section 4 recommends the use of (4.6) while the present section recommends against use of (6.6). At the risk of beating a dead horse, the reason is that section 4 assumed that models are non-nested, while the present section assumes that they are nested.
7. The restriction to one step ahead forecasts is for the following reason. For multiple step forecasts, the difference between model 1 and model 2 MSPEs presumably has a negative expectation. And simulations in Clark and McCracken (2003) generally find that use of standard critical values results in too few rejections. But sometimes there are too many rejections. This apparently results because of problems with HAC estimation of the standard error of the MSPE difference (private communication from Todd Clark).
8. Note, however, that Clark and McCracken (2001, p90) report simulation evidence that this statistic has less power than the non-normal statistics discussed above.
9. The reader may wonder whether asymptotic normality violates the rule of thumb enunciated at the beginning of this section, because $f_t = e_{1t}' X_{2t}' \beta_2^*$ is identically zero when evaluated at population $\beta_2^* = 0$. At the risk of confusing rather than clarifying, let me briefly note that the rule of thumb still applies, but only with a twist on the conditions given in the previous section. This twist, which is due to Giacomini and White (2003), holds R fixed as the sample size grows. Thus in population the random variable of interest is $f_t = e_{1t}' X_{2t}' \hat{\beta}_{2t}$, which for the fixed or rolling schemes is nondegenerate for all t . (Under the recursive scheme, $\hat{\beta}_{2t} \rightarrow_p 0$ as $t \rightarrow \infty$, which implies that f_t is degenerate for large t .) It is to be emphasized that technical conditions (R fixed vs. $R \rightarrow \infty$) are not arbitrary. Reasonable technical conditions should reasonably rationalize finite sample behavior. For tests of equal MSPE and encompassing discussed in the first part of this section, a vast range of simulation evidence suggests that the $R \rightarrow \infty$ condition generates a reasonably accurate asymptotic approximation (i.e., non-normality is implied by the theory and is found in the simulations.) The more modest array of simulation evidence for the R fixed approximation suggests that this approximation might work tolerably for the moment $E e_{1t}' X_{2t}' \beta_2^*$, provided the rolling or fixed scheme is used.

REFERENCES

Andrews, Donald W.K., 1991, "Heteroskedasticity and Autocorrelation Consistent Covariance Matrix Estimation," *Econometrica* 59, 1465-1471.

Andrews, Donald W.K. and J. Christopher Monahan, 1991, "An Improved Heteroskedasticity and Autocorrelation Consistent Covariance Matrix Estimator," *Econometrica* 60, 953-66.

Ashley. R., Granger, Clive .W.J. and Richard Schmalensee, 1980, "Advertising and Aggregate Consumption: An Analysis of Causality," *Econometrica* 48, 1149-1168.

Avramov, Doron, 2002, "Stock Return Predictability and Model Uncertainty", *Journal of Financial Economics* 64, 423-458.

Chao, John, Valentina Corradi and Norman R. Swanson, 2001, "Out-Of-Sample Tests for Granger Causality," *Macroeconomic Dynamics* 5, 598-620.

Chen, Shiu-Sheng, 2004, "A Note on In-Sample and Out-of-Sample Tests for Granger Causality," forthcoming, *Journal of Forecasting*.

Cheung, Yin-Wong, Menzie D. Chinn and Antonio Garcia Pascual, 2003, "Empirical Exchange Rate Models of the Nineties: Are Any Fit to Survive?", forthcoming, *Journal of International Money and Finance*.

Chong, Y.Y. and David F. Hendry, 1986, "Econometric evaluation of linear macro-economic models", *Review of Economic Studies*, 53, 671-690.

Christiano, Lawrence J., 1989, "P*: Not the Inflation Forecaster's Holy Grail," *Federal Reserve Bank of Minneapolis Quarterly Review* 13, 3-18.

Clark, Todd E. and Michael W. McCracken, 2001, "Tests of Equal Forecast Accuracy and Encompassing for Nested Models," *Journal of Econometrics* 105, 85-110.

Clark, Todd E. and Michael W. McCracken, 2003, "Evaluating Long Horizon Forecasts," manuscript, University of Missouri.

Clark, Todd E. and Michael W. McCracken, 2004, "Improving Forecast Accuracy by Combining Recursive and Rolling Forecasts," manuscript, University of Missouri.

Clark, Todd E. and Michael W. McCracken, 2005, "The Power of Tests of Predictive Ability in the Presence of Structural Breaks", *Journal of Econometrics*, 124, 1-31.

Clark, Todd E. and Kenneth D. West, 2004, "Using Out-of-Sample Mean Squared Prediction Errors to Test the Martingale Difference Hypothesis," forthcoming, *Journal of Econometrics*.

Clements Michael P and A.B. Galvao, 2004, "A Comparison of Tests of Nonlinear Cointegration with Application to the Predictability of Us Interest Rates Using the Term Structure," *International Journal of Forecasting* 20, 219-236.

- Corradi, Valentini and Norman R. Swanson, 2004a, "Bootstrap Procedures for Recursive Estimation Schemes With Applications to Forecast Model Selection," manuscript, Rutgers University.
- Corradi, Valentini and Norman R. Swanson, 2004b, "Predictive Density Evaluation," chapter in Handbook of Forecasting.
- Corradi, Valentini Norman R. Swanson and Claudia Olivetti, 2001, "Predictive Ability with Cointegrated Variables," *Journal of Econometrics* 104, 315-358.
- Davidson, Russell and James G. MacKinnon, 1984, "Model Specification Tests Based on Artificial Linear Regressions," *International Economic Review* 25 , 485-502.
- den Haan, Wouter J, and Andrew T. Levin, 2000, "Robust Covariance Matrix Estimation with Data-Dependent VAR Prewhitening Order," NBER Technical Working Paper: 255.
- Diebold, Francis X. and Robert S. Mariano, 1995, "Comparing Predictive Accuracy," *Journal of Business and Economic Statistics* 13, 253-263.
- Fair, Ray .C., 1980, "Estimating the predictive accuracy of econometric models," *International Economic Review*, 21, 355-378.
- Faust, Jon, John H. Rogers and Jonathan H. Wright, 2004, "News and Noise in G-7 GDP Announcements," forthcoming, *Journal of Money, Credit and Banking*.
- Ferreira, Miguel A., 2004, "Forecasting the Comovements of Spot Interest Rates," forthcoming, *Journal of International Money and Finance*.
- Ghysels, Eric and Alastair Hall, 1990, "A Test for Structural Stability of Euler Conditions Parameters Estimated via the Generalized Method of Moments Estimator," *International Economic Review* 31, 355-364.
- Giacomini, Raffaella and Halbert White, 2003, "Tests of Conditional Predictive Ability," manuscript, University of California at San Diego.
- Granger, C.W.J and Paul Newbold, 1977, *Forecasting Economic Time Series*, New York: Academic Press.
- Hansen, Lars Peter, 1982, "Large Sample Properties of Generalized Method of Moments Estimators," *Econometrica* 50, 1029-54.
- Hansen, Peter Reinhard, 2003, "A Test for Superior Predictive Ability," manuscript, Stanford University.
- Hansen, Peter Reinhard, Asger Lunde and James Nason, 2004, "Model Confidence Sets for Forecasting Models," manuscript, Stanford University.
- Harvey, David I., Stephen J. Leybourne, and Paul Newbold, 1998, "Tests for Forecast Encompassing," *Journal of Business and Economic Statistics* 16, 254-59.
- Hueng, C. James, and Ka Fu Wong , 2000, "Predictive Abilities of Inflation-Forecasting Models Using

Real Time Data,” Working Paper No 00-10-02 The University of Alabama.

Hoffman, Dennis L. and Adrian R. Pagan, 1989, “Practitioners Corner: Post Sample Prediction Tests for Generalized Method of Moments Estimators,” *Oxford Bulletin of Economics and Statistics* 51, 333-343.

Hong, Yongmiao and Tae-Hwy Lee, 2003, “Inference on Predictability of Foreign Exchange Rates via Generalized Spectrum and Nonlinear Time Series Models,” *Review of Economics and Statistics* 85, 1048-62.

Hueng, C. James, 1999, “Money Demand in an Open-Economy Shopping-Time Model: An Out-of-Sample-Prediction Application to Canada,” *Journal of Economics and Business* 51, 489-503.

Inoue, Atsushi, and Lutz Kilian, 2004a, “In-Sample or Out-of-Sample Tests of Predictability: Which One Should We Use?,” forthcoming, *Econometric Reviews*.

Inoue, Atsushi, and Lutz Kilian, 2004b, “On the Selection of Forecasting Models,” manuscript, University of Michigan.

Leitch, Gordon and J. Ernest Tanner, 1991, “Economic Forecast Evaluation: Profits versus the Conventional Error Measures,” *American Economic Review* 81, 580-590.

Lettau, Martin and Sydney Ludvigson, 2001, “Consumption, Aggregate Wealth, and Expected Stock Returns,” *Journal of Finance* 56, 815-849.

Mark, Nelson, 1995, “Exchange Rates and Fundamentals: Evidence on Long-Horizon Predictability,” *American Economic Review* 85, 201-218.

McCracken, Michael W., 2000, “Robust Out of Sample Inference,” *Journal of Econometrics* 99, 195-223.

McCracken, Michael W., 2004, “Asymptotics for Out of Sample Tests of Causality,” manuscript, University of Missouri.

McCracken, Michael W., and Stephen Sapp, 2003, “Evaluating the Predictability of Exchange Rates Using Long Horizon Regressions,” forthcoming, *Journal of Money, Credit and Banking*.

Meese, Richard A., and Kenneth Rogoff, 1983, “Empirical Exchange Rate Models of the Seventies: Do They Fit Out of Sample?” *Journal of International Economics* 14, 3-24.

Meese, Richard A., and Kenneth Rogoff, 1988, “Was it Real? The Exchange Rate - Interest Differential over the Modern Floating Rate Period,” *Journal of Finance* 43, 933-948.

Mizrach, Bruce, 1995, “Forecast Comparison in L_2 ,” manuscript, Rutgers University.

Morgan, W.A., 1939, “A test for significance of the difference between two variances in a sample from a normal bivariate population,” *Biometrika* 31, 13-19.

Newey, Whitney K. and Kenneth D. West, 1987, “A Simple, Positive Semidefinite, Heteroskedasticity and Autocorrelation Consistent Covariance Matrix,” *Econometrica* 55, 703-708.

- Newey, Whitney K. and Kenneth D. West, 1994, "Automatic Lag Selection in Covariance Matrix Estimation," *Review of Economic Studies* 61, 631-654.
- Pagan, Adrian R. and Anthony D. Hall, 1983, "Diagnostic Tests as Residual Analysis," *Econometric Reviews* 2, 159-218.
- Politis, D. N. and Joseph P. Romano, 1994. "The Stationary Bootstrap." *Journal of the American Statistical Association* 89, 1301-1313.
- Romano, Joseph P. and Micahel Wolf, 2003, "Stepwise Multiple Testing as Formalize Data Snooping," manuscript, Stanford University.
- Rossi, Barbara, 2003, "Testing Long-horizon Predictive Ability with High Persistence, and the Meese-Rogoff Puzzle," forthcoming *International Economic Review*.
- Sarno, Lucia, Daniel L. Thornton and Giorgio Valente, "Federal Funds Rate Prediction" forthcoming, *Journal of Money, Credit and Banking*.
- Shintani, Mototsugu , 2004, "Nonlinear Analysis of Business Cycles Using Diffusion Indexes: Applications to Japan and the US," forthcoming, *Journal of Money, Credit and Banking*.
- Stock, James H. and Mark W. Watson, 1999, "Forecasting Inflation," *Journal of Monetary Economics*, 44, 293-335.
- Stock, James H. and Mark W. Watson, 2002, "Macroeconomic Forecasting Using Diffusion Indexes," *Journal of Business and Economic Statistics* 20, 147-162.
- Storey, John D., 2002, "A Direct Approach to False Discovery Rates," *Journal of the Royal Statistical Society, Series B* 64 479-498.
- West, Kenneth D., 1996, "Asymptotic Inference About Predictive Ability," *Econometrica* 64 , 1067-1084.
- West, Kenneth D., 2001, "Tests of Forecast Encompassing When Forecasts Depend on Estimated Regression Parameters," *Journal of Business and Economic Statistics* 19, 29-33.
- West, Kenneth D. and Dongchul Cho, 1995, "The Predictive Ability of Several Models of Exchange Rate Volatility," *Journal of Econometrics* 69, 367-391.
- West, Kenneth D., Hali J. Edison and Dongchul Cho, 1993, "A Utility Based Comparison of Some Models of Exchange Rate Volatility," *Journal of International Economics* 35, 23-46.
- West, Kenneth D. and Michael W. McCracken, 1998, "Regression Based Tests of Predictive Ability," *International Economic Review* 39, 817-840.
- White, Halbert, 1984, "Asymptotic Theory for Econometricians," New York: Academic Press.
- White, Halbert, 2000, "A Reality Check for Data Snooping," *Econometrica* 68, 1097-1126.

Wilson, Edwin B., 1934, "The Periodogram of American Business Activity," *The Quarterly Journal of Economics* 48, 375-417.

Wooldridge, Jeffrey M., 1990, "A Unified Approach to Robust, Regression-Based Specification Tests," *Econometric Theory* 6, 17-43.

Table 1

Sample Analogues for λ_{fh} , λ_{hh} and λ

	Recursive	Rolling, $P \leq R$	Rolling, $P > R$	Fixed
λ_{fh}	$1 - \frac{R}{P} \ln(1 + \frac{P}{R})$	$\frac{1}{2} \frac{P}{R}$	$1 - \frac{1}{2} \frac{R}{P}$	0
λ_{hh}	$2[1 - \frac{R}{P} \ln(1 + \frac{P}{R})]$	$\frac{P}{R} - \frac{1}{3} \frac{P^2}{R^2}$	$1 - \frac{1}{3} \frac{R}{P}$	$\frac{P}{R}$
λ	1	$1 - \frac{1}{3} \frac{P^2}{R^2}$	$\frac{2R}{3P}$	$1 + \frac{P}{R}$

Notes:

1. The recursive, rolling and fixed schemes are defined in section 4 and illustrated for an AR(1) in equation(4.2).
2. P is the number of predictions, R the size of the smallest regression sample. See section 4 and equation (4.1).
3. The parameters λ_{fh} , λ_{hh} and λ are used to adjust the asymptotic variance covariance matrix for uncertainty about regression parameters used to make predictions. See section 5 and Tables 2 and 3.

Table 2

Recommended Sources of Critical Values, Small Number of Models

Source of critical values	Conditions for use
<p>A. Use critical values associated with asymptotic normality, abstracting from any dependence of predictions on estimated regression parameters, as illustrated for scalar hypothesis test in (3.2) and a vector test in (4.11).</p>	<ol style="list-style-type: none"> 1. Prediction sample size P is small relative to regression sample size R, say $P/R < 0.1$ (any sampling scheme or moment, nested or nonnested models). 2. MSPE or mean absolute error in nonnested models. 3. Sampling scheme is recursive, moment of interest is mean prediction error or correlation between a given model's prediction error and prediction. 4. Sampling scheme is recursive, one step ahead conditionally homoskedastic prediction errors, moment of interest is either: (a) first order autocorrelation or (b) encompassing in the form (4.7c). 5. MSPE, one step ahead forecasts, nested models, equality of MSPE rejects (implying that it will also reject with an even smaller p-value if an asymptotically valid test is used).
<p>B. Use critical values associated with asymptotic normality, but adjust test statistics to account for the effects of uncertainty about regression parameters.</p>	<ol style="list-style-type: none"> 1. Mean prediction error, first order autocorrelation of one step ahead prediction errors, zero correlation between a prediction error and prediction, encompassing in the form (4.7c) (with the exception of point C3), encompassing in the form (4.7d) for nonnested models. 2. Zero correlation between a prediction error and another model's vector of predictors (nested or nonnested) (Chao et al. (2001)). 3. A general vector of moments or a loss or utility function that satisfies a suitable rank condition. 4. Rolling or fixed scheme, MSPE, nested models with the null model a random walk: adjust the point estimates and standard errors in a certain way (Clark and West (2004)).
<p>C. Use non-standard critical values.</p>	<ol style="list-style-type: none"> 1. MSPE or encompassing in the form (4.7d), in nested models, one step ahead prediction errors: use critical values from McCracken (2004) or Clark and McCracken (2001). 2. MSPE, encompassing in the form (4.7d) or mean absolute error, in nested models, and in contexts not covered by A5, B4 or C1 (e.g., multistep predictions with the null model not a random walk): simulate/bootstrap your own critical values. 3. Recursive scheme, $\beta_1^*=0$, encompassing in the form (4.7c): simulate/bootstrap your own critical values.

Note: Rows B and C assume that P/R is sufficiently large, say $P/R \geq 0.1$, that there may be nonnegligible effects of estimation uncertainty about parameters used to make forecasts. The results in Row A, points 2 through 5, apply whether or not P/R is large.

Table 3

Recommended Procedures, Small Number of Models

A. Tests of Adequacy of a Single Model, $y_t = X_t' \beta^* + e_t$

(1) Description	(2) Null hypothesis	(3) Recommended procedure	(4) Asymptotic normal critical values?
1. Mean prediction error (bias)	$E(y_t - X_t' \beta^*) = 0$, or $Ee_t = 0$	Regress prediction error on a constant, divide HAC t-stat by $\sqrt{\lambda}$.	Y
2. Correlation between prediction error and prediction (efficiency)	$E(y_t - X_t' \beta^*) X_t' \beta^* = 0$, or $Ee_t X_t' \beta^* = 0$	Regress \hat{e}_{t+1} on $X_{t+1}' \hat{\beta}_t$, divide HAC t-stat by $\sqrt{\lambda}$, or regress y_{t+1} on prediction $X_{t+1}' \hat{\beta}_t$, divide HAC t-stat (for testing coefficient value of 1) by $\sqrt{\lambda}$.	Y
3. First order correlation of one step ahead prediction errors	$E(y_{t+1} - X_{t+1}' \beta^*)(y_t - X_t' \beta^*) = 0$, or $Ee_{t+1} e_t = 0$.	a. Prediction error conditionally homoskedastic: <ol style="list-style-type: none"> 1. Recursive scheme: regress \hat{e}_{t+1} on \hat{e}_t, use OLS t-stat. 2. Rolling or fixed schemes: regress \hat{e}_{t+1} on \hat{e}_t and X_t, use OLS t-stat on coefficient on \hat{e}_t. b. Prediction error conditionally heteroskedastic: adjust standard errors as described in section 5 above.	Y

Notes:

1. The quantity λ is computed as described in Table 1. ‘‘HAC’’ refers to a heteroskedasticity and autocorrelation consistent covariance matrix. Throughout, it is assumed that predictions rely on estimated regression parameters and that P/R is large enough, say $P/R \geq 0.1$, that there may be nonnegligible effects of such estimation. If P/R is small, say $P/R < 0.1$, any such effects may well be negligible, and one can use standard results as described in sections 3 and 4.

B. Tests Comparing a Pair of Nonnested Models, $y_i = X_{1t}'\beta_1^* + e_{1t}$ vs. $y_i = X_{2t}'\beta_2^* + e_{2t}$, $X_{1t}'\beta_1^* \neq X_{2t}'\beta_2^*$, $\beta_2^* \neq 0$

(1)

(2)

(3)

(4)

Description

Null hypothesis

Recommended procedure

Asymptotic normal critical values?

1. Mean squared prediction error (MSPE)	$E(y_i - X_{1t}'\beta_1^*)^2 - E(y_i - X_{2t}'\beta_2^*)^2 = 0$, or $Ee_{1t}^2 - Ee_{2t}^2 = 0$	Regress $\hat{e}_{1t+1}^2 - \hat{e}_{2t+1}^2$ on a constant, use HAC t-stat.	Y
2. Mean absolute prediction error (MAPE)	$E y_i - X_{1t}'\beta_1^* - E y_i - X_{2t}'\beta_2^* = 0$, or $E e_{1t} - E e_{2t} = 0$	Regress $ \hat{e}_{1t} - \hat{e}_{2t} $ on a constant, use HAC t-stat.	Y
3. Zero correlation between model 1's prediction error and the prediction from model 2 (forecast encompassing)	$E(y_i - X_{1t}'\beta_1^*)X_{2t}'\beta_2^* = 0$, or $Ee_{1t}X_{2t}'\beta_2^* = 0$	a. Recursive scheme, prediction error e_{1t} , homoskedastic conditional on both X_{1t} and X_{2t} : regress \hat{e}_{1t+1} on $X_{2t+1}'\hat{\beta}_{2t}$, use OLS t-stat. b. Recursive scheme, prediction error e_{1t} conditionally heteroskedastic, or rolling or fixed scheme: regress \hat{e}_{1t+1} on $X_{2t+1}'\hat{\beta}_{2t}$ and X_{1t} , use HAC t-stat on coefficient on $X_{2t+1}'\hat{\beta}_{2t}$.	Y
4. Zero correlation between model 1's prediction error and the difference between the prediction errors of the two models (another form of forecast encompassing)	$E(y_i - X_{1t}'\beta_1^*) \times [(y_i - X_{1t}'\beta_1^*) - (y_i - X_{2t}'\beta_2^*)] = 0$, or $Ee_{1t}(e_{1t} - e_{2t}) = 0$	Adjust standard errors as described in section 5 above and illustrated in West (2001).	Y
5. Zero correlation between model 1's prediction error and the model 2 predictors	$E(y_i - X_{1t}'\beta_1^*)X_{2t} = 0$, or $Ee_{1t}X_{2t} = 0$	Adjust standard errors as described in section 5 above and illustrated in Chao et al. (2001).	Y

See notes to Table 3A.

C. Tests of Comparing a Pair of Nested Models, $y_t = X_{1t}'\beta_1^* + e_{1t}$ vs. $y_t = X_{2t}'\beta_2^* + e_{2t}$, $X_{1t} \subset X_{2t}$, $X_{2t}' = (X_{1t}', X_{22t}')'$

(1) (2) (3) (4)

Description Null hypothesis Recommended procedure Asymptotic normal critical values?

1. Mean squared prediction error (MSPE)	$E(y_t - X_{1t}'\beta_1^*)^2 - E(y_t - X_{2t}'\beta_2^*)^2 = 0$, or $Ee_{1t}^2 - Ee_{2t}^2 = 0$	a. One step ahead, conditionally homoskedastic disturbances: use critical values from McCracken (2004).	N
		b. One step ahead, equality of MSPE rejects (implying that it will also reject with an even smaller p-value if an asymptotically valid test is used).	Y
		c. Rolling or fixed scheme, null model is a random walk: adjust point estimate and standard error in a certain way (Clark and West (2004)).	Y
		d. Simulate/bootstrap your own critical values.	N
2. Mean absolute prediction error (MAPE)	$E y_t - X_{1t}'\beta_1^* - E y_t - X_{2t}'\beta_2^* = 0$, or $E e_{1t} - E e_{2t} = 0$	Simulate/bootstrap your own critical values.	N
3. Zero correlation between model 1's prediction error and the prediction from model 2 (forecast encompassing)	$E(y_t - X_{1t}'\beta_1^*)X_{2t}'\beta_2^* = 0$, or $Ee_{1t}X_{2t}'\beta_2^* = 0$	a. $\beta_1^* \neq 0$: regress \hat{e}_{1t+1} on $X_{2t+1}'\hat{\beta}_{2t}$, divide HAC t-stat by $\sqrt{\lambda}$.	Y
		b. $\beta_1^* = 0$ ($\Rightarrow \beta_2^* = 0$), 1. Rolling or fixed scheme: regress \hat{e}_{1t+1} on $X_{2t+1}'\hat{\beta}_{2t}$, use HAC t-stat.	Y
		2. $\beta_1^* = 0$, recursive scheme: simulate/bootstrap your own critical values.	N
4. Zero correlation between model 1's prediction error and the difference between the prediction errors of the two models (another form of forecast encompassing)	$E(y_t - X_{1t}'\beta_1^*) \times [(y_t - X_{1t}'\beta_1^*) - (y_t - X_{2t}'\beta_2^*)] = 0$ or $Ee_{1t}(e_{1t} - e_{2t}) = 0$	a. One step ahead, conditionally homoskedastic disturbances: use critical values from Clark and McCracken (2001).	N
		b. Simulate/bootstrap your own critical values.	N
5. Zero correlation between model 1's prediction error and the model 2 predictors	$E(y_t - X_{1t}'\beta_1^*)X_{22t} = 0$, or $Ee_{1t}X_{22t} = 0$	Adjust standard errors as described in section 5 above and illustrated in Chao et al. (2001).	Y

1. See notes to Table 3A.

2. Under the null, the coefficients on X_{22t} (the regressors included in model 2 but not model 1) are zero. Thus, $X_{1t}'\beta_1^* = X_{2t}'\beta_2^*$ and $e_{1t} = e_{2t}$.