



ECONOMIC RESEARCH
FEDERAL RESERVE BANK OF ST. LOUIS
WORKING PAPER SERIES

Tests of Conditional Predictive Ability: Some Simulation Evidence

Authors	Michael W. McCracken
Working Paper Number	2019-011C
Revision Date	April 2019
Citable Link	https://doi.org/10.20955/wp.2019.011
Suggested Citation	McCracken, M.W., 2019; Tests of Conditional Predictive Ability: Some Simulation Evidence, Federal Reserve Bank of St. Louis Working Paper 2019-011. URL https://doi.org/10.20955/wp.2019.011

Federal Reserve Bank of St. Louis, Research Division, P.O. Box 442, St. Louis, MO 63166

The views expressed in this paper are those of the author(s) and do not necessarily reflect the views of the Federal Reserve System, the Board of Governors, or the regional Federal Reserve Banks. Federal Reserve Bank of St. Louis Working Papers are preliminary materials circulated to stimulate discussion and critical comment.

Tests of Conditional Predictive Ability: Some Simulation Evidence *

Michael W. McCracken
Federal Reserve Bank of St. Louis

May 2019

Abstract

In this note we use simple examples and associated simulations to investigate the size and power properties of tests of predictive ability described in Giacomini and White (2006; *Econometrica*). While we find that the tests can be accurately sized and powerful in large enough samples we identify details associated with the tests that are not otherwise apparent from the original text. In order of importance these include (i) the proposed test of equal finite-sample unconditional predictive ability is not asymptotically valid under the fixed scheme, (ii) for the same test, but when the rolling scheme is used, very large bandwidths are sometimes required when estimating long-run variances, and (iii) when conducting the proposed test of equal finite-sample conditional predictive ability, conditional heteroskedasticity is likely present when lagged loss differentials are used as instruments.

JEL Nos.: C53, C12, C52

Keywords: prediction, out-of-sample, inference

* *McCracken*: Research Division; Federal Reserve Bank of St. Louis; P.O. Box 442; St. Louis, MO 63166; michael.w.mccracken@stls.frb.org. The views expressed herein are solely those of the authors and do not necessarily reflect the views of the Federal Reserve Bank of St. Louis or the Federal Reserve System.

1 Introduction

It has become commonplace to evaluate the quality of a predictive model based upon its ability to forecast in a pseudo-out-of-sample framework. This approach often consists of splitting the existing time series of data into two portions: an in-sample portion used to estimate model parameters and a second portion used to construct and evaluate the accuracy of forecasts. In the context of point forecasts, the evaluation stage typically consists of (i) constructing forecast errors, (ii) forming sample averages of functions of these errors, and (iii) conducting inference related to the quality of the forecasts using said sample average.

An important step forward in this literature was made by Diebold and Mariano (1995). In the context of comparing the accuracy of two point forecasts, they show that a t-statistic associated with sample averages of loss differentials from these two models can be asymptotically Normal. Asymptotically valid inference on the mean of the loss differentials is therefore straightforward. This type of statistic was extended by West (1996) to explicitly account for the effect parameter estimation error has on the asymptotic variance of the average loss differential when estimated models are used to form the forecasts. This line of research then expanded to a variety of distinct topics including multiple testing (White 2000), comparisons between nested models (Clark and McCracken 2001, McCracken 2007), cointegrating relationships (Corradi, Olivetti, and Swanson, 2001), forecast breakdowns (Giacomini and Rossi, 2009), generated predictors (Goncalves, McCracken, and Perron 2017), generated predictands (Li and Patton, 2018), and many others.

In this literature, when the predictive models use estimated parameters, it is typically the case that the asymptotic theory assumes that the sample sizes used to estimate the model parameters are large enough that the parameter estimates ($\hat{\beta}_t$) are consistent for their population counterparts (β^*). In practice this means test statistics, that are functions of the finite-sample h -step ahead forecast errors $u_{t+h}(\hat{\beta}_t) = \hat{u}_{t+h}$, are being used to test null hypotheses related to population-level forecast errors $u_{t+h}(\beta^*) = u_{t+h}$.

A different approach is taken by Giacomini and White (2006, henceforth GW). They emphasize the need to evaluate the accuracy of forecasts as we observe them in practice which requires admitting that they typically depend on estimated parameters. As an example, when comparing the accuracy of point forecasts from two models ($i = 1, 2$) under quadratic loss, the bulk of the literature focuses on a null hypothesis of the form $E(u_{1,t+h}^2 - u_{2,t+h}^2) = 0$ while GW instead consider hypotheses of the form $E(\hat{u}_{1,t+h}^2 - \hat{u}_{2,t+h}^2) = 0$ and $E(\hat{u}_{1,t+h}^2 - \hat{u}_{2,t+h}^2 | \mathfrak{F}_t) = 0$. The first of these three hypotheses

is a statement about the unconditional expectation of the population-level loss differential while the second and third are statements about the unconditional and conditional (based on a time t information set \mathfrak{S}_t) expectation of the finite-sample loss differential.¹

These two new hypotheses required new test statistics and more importantly, a new approach to inference that prevents the parameter estimates from converging to their population counterparts as the total sample size increases. GW achieve this by requiring that a finite number of observations (R) are always used to estimate the parameters.² Either the estimates depend on a rolling window of observations x such that $\hat{\beta}_t = \beta(x_{t-R+1}, \dots, x_t)$ or a fixed window of observations such that $\hat{\beta}_t = \beta(x_1, \dots, x_R)$. With this assumption in hand, and along with other technical assumptions on the loss differential $\hat{d}_{t+h} = \hat{u}_{1,t+h}^2 - \hat{u}_{2,t+h}^2$, $t = R, \dots, T - h = R + P - 1$, they prove that a statistic of the form $P^{-1/2} \sum_{t=R}^{T-h} \hat{d}_{t+h} / \hat{\omega}$ is asymptotically standard Normal under the null hypothesis of equal finite sample predictive ability (either unconditional $E(\hat{u}_{1,t+h}^2 - \hat{u}_{2,t+h}^2) = 0$ or conditional $E(\hat{u}_{1,t+h}^2 - \hat{u}_{2,t+h}^2 | \mathfrak{S}_t) = 0$) so long as a consistent estimate of the long-run variance of \hat{d}_{t+h} , ω^2 , is used to form $\hat{\omega}$.³

In this paper we investigate the size and power properties of the tests of equal finite sample predictive ability using simulations of data generating processes (DGPs) that can be shown to satisfy the null hypotheses. Our simulation are sometimes, but not always, supportive of the size and power properties reported in GW.⁴ When the forecast horizon is one, the nominal size of the test can be good for tests of conditional, finite-sample predictive ability so long as conditional heteroskedasticity is properly accounted for. As the horizon grows, size distortions arise due to the need to estimate long-run variances. In accordance with the theory, these size distortions dissipate as the sample size increases.

Size distortions are particularly acute when we consider the tests of unconditional, finite-sample predictive ability. For these tests there are two issues that cause size distortions. The first is that under the fixed scheme, the test is not asymptotically Normal even under the null because the loss differentials are not asymptotically independent. In fact, the

¹As described in equation (2) on page 1549 of GW.

²As a technical matter, they permit R to vary across forecast origins but the maximum is bounded from above. For brevity we assume it is constant.

³In GW, the hypotheses are stated more broadly and apply to any well-behaved loss function and are not exclusive to point forecasts. Even so, for exposition purposes we will focus on the most common applications in which quadratic loss is used to evaluate point forecasts.

⁴The simulations in GW are not based on DGPs that satisfy either of the null hypotheses. Their first simulation satisfies the null of equal unconditional predictive ability *on average* rather than for all t . Their second simulation does not include a DGP that implies loss differentials that satisfy the null of equal conditional predictive ability.

statistic diverges with probability one under the null. This is despite the fact that the loss differentials are covariance stationary. As such we rebut the claim made in Diebold (2015) that so long as loss differentials are covariance stationary, the t-statistic delineated in Diebold and Mariano (1995) *must* be asymptotically Normal. The issue is not whether or not the loss differentials are precisely covariance stationary but whether or not the loss differentials are asymptotically independent due to properties related to mixing, near-epoch dependence, etc.

The second reason for the size distortions is that under the rolling scheme, estimating long-run variances is particularly challenging. This arises due to the fact that even one step ahead forecasts imply loss differentials that are serially correlated of order $R - 1$. Rolling window estimation of the parameters induces long-lagged serial correlation in the loss differentials and properly accounting for that degree of autocorrelation can be difficult in finite samples.

We also investigate the power of each test. To do so we emphasize the fact that, in the GW framework, satisfying the null hypothesis requires delineating DGPs that depend explicitly on the chosen value of R and the chosen forecasting scheme (fixed or rolling). We therefore consider two types of alternatives. The first is that the window size used to estimate the parameters (\tilde{R}) does not align with that maintained in the DGP (R). For the second, we consider the case in which the fixed (rolling) scheme is used to estimate the model parameters while the DGP is defined using the rolling (fixed) scheme. For the former, reasonable power only arises when the out-of-sample period is quite large and even then it is rarely the case that the actual power is greater than 70%. In many cases, actual power of the test aligns closely with rejection frequencies observed under the null. For the latter alternative, actual power is sometimes improved and sometimes not. Regardless, rejection frequencies greater than 70% are generally restricted to out-of-sample sizes that are quite large.

Before proceeding it is worth noting others who have contributed to the literature on tests designed to evaluate whether two models exhibit equal finite-sample, rather than population predictive ability. In the context of developing a test of pairwise forecast comparisons in an unstable environment, Giacomini and Rossi (2010) delineate a DGP that satisfies the null of equal unconditional finite-sample predictive ability $E(\hat{u}_{1,t+1}^2 - \hat{u}_{2,t+1}^2) = 0$ under the rolling scheme. They do not provide comparable evidence for the fixed scheme. Clark and McCracken (2015) develop a test of equal average unconditional finite-sample predictive

ability $E[P^{-1} \sum_{t=R}^{T-h} (\hat{u}_{1,t+h}^2 - \hat{u}_{2,t+h}^2)] = 0$ for both rolling and expanding estimation windows. Simulation evidence is provided using DGPs but the null is distinct from that in GW. Coroneo and Iacone (2018) investigate fixed-smoothing approaches to inference when testing the null of equal unconditional finite-sample predictive ability $E(\hat{u}_{1,t+1}^2 - \hat{u}_{2,t+1}^2) = 0$ but do so without specifying DGPs that imply loss differentials that satisfy the null. Perhaps most importantly, Timmermann and Zhu (2017) develop theoretical results allowing expanding estimation windows when testing the same null hypothesis $E(\hat{u}_{1,t+h}^2 - \hat{u}_{2,t+h}^2 | \mathfrak{F}_t) = 0$ considered in GW. Despite their theoretical results, they do not provide simulations that explicitly delineate DGPs that satisfy the null hypothesis.⁵

The remainder proceeds as follows. Section 2 provides simulation evidence on the size and power of tests of the null hypothesis $E(\hat{u}_{1,t+h}^2 - \hat{u}_{2,t+h}^2) = 0$ for all $t = R, \dots, T - h$ based on a fully specified DGP. Section 3 does the same but for the null hypothesis $E(\hat{u}_{1,t+h}^2 - \hat{u}_{2,t+h}^2 | \mathfrak{F}_t) = 0$. Section 4 concludes.

2 Tests of Unconditional, Finite-Sample Predictive Ability

In this section we delineate a DGP that satisfies the null hypothesis $E(\hat{u}_{1,t+h}^2 - \hat{u}_{2,t+h}^2) = 0$ for all $t = R, \dots, T - h$. We do so under both the fixed and rolling estimation scheme and across a range of forecast horizons h , sample sizes T , and window sizes R . By varying the parameterization of the DGP we also address the power properties of the test. Specifically, in this section we consider a number of permutations in which the estimation window \tilde{R} does not equal the value of R that defines the DGP under the null. In each of these instances the null is not satisfied and constitutes an alternative.

- Our DGP is motivated by simple application in which a no-change forecast is compared to that from a location model with a single estimated parameter. That is, $\hat{y}_{1,t+h} = 0$ while $\hat{y}_{2,t+h} = \bar{y}_t$ where $\bar{y}_t = \bar{y}_{\tilde{R}}$ for all t under the fixed scheme and $\bar{y}_t = \tilde{R}^{-1} \sum_{s=t-\tilde{R}+1}^t y_s$ under the rolling scheme. Straightforward algebra reveals that if $y_t = \mu + \eta_t$, $\eta_t = \varepsilon_t + \sum_{j=1}^{h-1} \theta_j \varepsilon_{t-j}$ with $\varepsilon_t \sim i.i.d.N(0, \sigma^2)$ then $E(\hat{u}_{1,t+h}^2 - \hat{u}_{2,t+h}^2) = 0$ for all $t = R, \dots, T - h$ when $\mu = R^{-1/2}[\gamma_0 + 2 \sum_{j=1}^{R-1} \left(\frac{R-j}{R}\right) \gamma_j]^{1/2}$, $\gamma_j = E\eta_t \eta_{t-j}$, and $\tilde{R} = R$.
- In all simulations the test statistic takes the form $P^{-1/2} \sum_{t=\tilde{R}}^{T-h} \hat{d}_{t+h} / \hat{\omega}$. With this

⁵Timmermann and Zhu (2017) describe a simple DGP involving two non-nested models in the text of their paper but do not consider this example in their simulations. Since the described forecasting exercise does not entail estimating model parameters, a focus of this note, we do not pursue this DGP in our simulations.

statistic we use standard Normal critical values to conduct a two-sided test of the null. In all tables we only report results associated with a nominal size of 5%. Unreported results associated with 1% and 10% provide comparable results and are excluded for brevity.

- We use sample sizes in which R , \tilde{R} , and P range across 25, 75, 25, and 175 – similar to those used in GW. We also consider $P = 1000$. In each simulation, the total sample size T is $\tilde{R} + P + h - 1$.
- Throughout we set the variance of the primitive shocks σ^2 to 1. Unreported results that set σ^2 to 0.1 and 10 are quite similar and excluded for brevity. The MA coefficients are set to $\theta_j = (0.5)^j$ in all simulations.
- We consider forecast horizons h of 1, 3, and 12. Unreported results that set h to 24 are comparable to those for 12 and excluded for brevity.
- In all results we use the Bartlett kernel to estimate ω^2 (i.e. Newey and West, 1987). The choice of bandwidth does have an affect and hence we experimented with a few including fixed values of h and \tilde{R} but also a sample dependent version $\lfloor 4(P/100)^{2/9} \rfloor + 1$ (Newey and West, 1994; Andrews and Monahan, 1992). With $h - 1$ lags, the actual size was better (than using $\tilde{R} - 1$) for smaller samples because it entailed estimating fewer autocovariances but was poor for the larger samples because it didn't estimate enough of them. With $\tilde{R} - 1$ lags, it was the reverse: actual size was worse (than using $h - 1$) in smaller samples but improved as the sample got larger because it was estimating the correct numbers of autocovariances.⁶ Broadly speaking, all size results were (weakly) more accurate using the sample dependent version and hence all reported results use that bandwidth.
- In all experiments the number of replications is 5,000.

2.1 Fixed

We begin by investigating size and power of the test when the fixed scheme is used. In Table 1, rows that are in bold are those for which the null hypothesis holds. All other values represent alternative hypotheses driven by the fact that $R \neq \tilde{R}$. In all cases we find that the actual size of the test is much higher than the nominal 5% level. Making it worse

⁶Note too that when $\tilde{R} > P$, a bandwidth of \tilde{R} observations is infeasible because there are no observations available to estimate, for example, the $\tilde{R} - 1$ th autocovariance.

is the fact that the rejection frequencies generally get worse, not better, as the sample size increases.

Given our simple DGPs, a bit of algebra allows us to identify the root of the size distortions. For simplicity let $h = 1$, $R = \tilde{R}$, and $\hat{\omega}^2 = P^{-1} \sum_{t=R}^{T-1} \hat{d}_{t+1}^2$. Rearranging terms we find that $P^{-1/2} \sum_{t=R}^{T-1} \hat{d}_{t+1} / \hat{\omega}$ equals

$$\frac{2(P^{-1/2} \sum_{t=R}^{T-1} \varepsilon_{t+1}) \bar{y}_R + P^{1/2} (\mu^2 - (R^{-1} \sum_{s=1}^R \varepsilon_s)^2)}{\sqrt{4\bar{y}_R^2 (P^{-1} \sum_{t=R}^{T-1} \varepsilon_{t+1}^2) + 4\bar{y}_R^2 (2\mu - \bar{y}_R) (P^{-1} \sum_{t=R}^{T-1} \varepsilon_{t+1}) + (\mu^2 - (R^{-1} \sum_{s=1}^R \varepsilon_s)^2)^2}}.$$

As P diverges, holding R constant, the denominator is $O_p(1)$ with probability limit $\sqrt{4\bar{y}_R^2 \sigma^2 + (\mu^2 - (R^{-1} \sum_{s=1}^R \varepsilon_s)^2)^2}$. In contrast, while the first part of the numerator is asymptotically Normal, the second part diverges to $\pm\infty$ with probability one and hence the test will always reject the null for large enough P .

To understand why the test statistic fails to be asymptotically Normal it is instructive to calculate the first and second moments of the loss differential. Given the definition of μ and the *i.i.d.* nature of ε_{t+1} , straightforward algebra reveals that $E\hat{d}_{t+1} = 0$, $E\hat{d}_{t+1}^2 = 4\sigma^2 E\bar{y}_R^2 + E(\mu^2 - (R^{-1} \sum_{s=1}^R \varepsilon_s)^2)^2$, and $E\hat{d}_{t+1} \hat{d}_{t+1-j} = E(\mu^2 - (R^{-1} \sum_{s=1}^R \varepsilon_s)^2)^2$ for all $j \geq 1$. Evidently, while the loss differential is covariance stationary, the autocovariances do not vanish as the lag increases. Since this implies that the loss differentials are not asymptotically independent, the average loss differential is not asymptotically Normal.

This example directly contradicts a point made in GW that the fixed scheme is compatible with their mixing conditions which in turn play a critical role in establishing asymptotic Normality of the average loss differential in their Theorem 4. While we can only speculate, it appears that GW misinterpret the role of the parameter τ in Lemma 2.1 of White and Domowitz (1984). Recall that White and Domowitz prove that for a mixing sequence Z_t , and finite positive integer τ , a finite dimensioned function $X_t = g(Z_t, \dots, Z_{t-\tau})$ is also mixing of the same order as Z_t . While it is true that under the fixed scheme, the loss differential \hat{d}_{t+1} is a finite dimensioned function of mixing variables (i.e. $\hat{d}_{t+1} = g(y_{t+1}, y_R, \dots, y_1)$ has $R+1$ arguments), the window size $(t+1) - 1$ is not finite as P increases and hence Lemma 2.1 of White and Domowitz (1984) doesn't apply. As such, Theorem 4 in GW is not, as currently proven, valid when using the fixed scheme.

As noted in the introduction, one can interpret these results as a proof-by-example negation of a claim made in Diebold (2015) that, so long as the loss differentials are covariance stationary, the t-statistic delineated in Diebold and Mariano (1995) *must* be asymptotically Normal. Although we agree that the loss differentials must have well-behaved moments for

the statistic to be asymptotically Normal, the loss differentials also need to be asymptotically independent. By its nature, using the fixed scheme imbeds an inherent lack of asymptotic independence into the loss differentials. As such it seems ill-advised to use the t-statistic advocated by Diebold and Mariano (1995) when the fixed scheme is used and the null hypothesis of interest is $E\hat{d}_{t+h} = 0$.

Returning now to the simulation results, since the actual size of the test is so poor, it is difficult to interpret actual power of the test when $R \neq \tilde{R}$. For the samples in which $P < 1000$, the rejection frequencies are often similar to those under the null unless $\tilde{R} = 25$ but even there actual power tops off around 40%. The rejection frequencies under the alternative rise to 70% and more when $P = 1000$ but the size distortions are increasing as well making interpretation difficult.

2.2 Rolling

In Table 2 we report similar results but when the rolling scheme is used. Here we find that the actual size of the test is much more in line with what we would expect from a nominally 5% test. At the shortest forecast horizon, $h = 1$, rejection frequencies are as high as 9% when $P = 25$ but otherwise range between 3% and 7%. At the longer horizons, the actual size of the test is as high as 20% but improve substantially as P increases.

The size distortions are due in part to the need to estimate a long-run variance that accounts for serial correlation in the loss differential. This serial correlation takes the form of an $MA(\tilde{R} - 1)$. To see this, consider the numerator of the test statistic when $R = \tilde{R}$, and $h = 1$. Rearranging terms we find that $P^{-1/2} \sum_{t=R}^{T-1} \hat{d}_{t+1}$ equals

$$2(P^{-1/2} \sum_{t=R}^{T-1} \varepsilon_{t+1} \bar{y}_t) + P^{-1/2} \sum_{t=R}^{T-1} (\mu^2 - (R^{-1} \sum_{s=t-R+1}^t \varepsilon_s)^2).$$

Both right hand side terms are zero mean and exhibit asymptotic independence and hence the numerator is asymptotically Normal. However, a closer look reveals that \hat{d}_{t+1} is much more persistent than one might have expected. The rolling windows approach induces serial correlation of order $\tilde{R} - 1$ which is substantially more persistent than one might have guessed given that the forecasts are at the one-step horizon. This, in turn, entails estimating a large number of autocovariances when forming a consistent estimate of the long-run variance.

Under the alternative in which $\tilde{R} \neq R$, power is quite poor with rejection frequencies generally below 20% for all but the largest sample size. In a few cases the test is biased:

rejection frequencies under the alternative are lower than under the null. For the largest sample size, actual power can be as large as 80% but remains as low as 10% for some permutations of $\tilde{R} \neq R$.

3 Tests of Conditional, Finite-Sample Predictive Ability

In this section we delineate a DGP that satisfies the null hypothesis $E(\hat{u}_{1,t+h}^2 - \hat{u}_{2,t+h}^2 | \mathfrak{S}_t) = 0$. In contrast to the previous section, the DGP is specific to whether the fixed or rolling scheme is being used. This means that if the DGP is designed relative to the fixed (rolling) scheme, but a rolling (fixed) model is used for forecasting, there must be a deviation from the null hypothesis. In addition, since $R \neq \tilde{R}$ continues to indicate a deviation from the null hypothesis, we provide simulations on power of the test against these two forms of alternatives.

Our DGPs are again developed based on a simple application in which a no-change forecast is compared to that from a location model with a single estimated parameter. That is, $\hat{y}_{1,t+h} = 0$ while $\hat{y}_{2,t+h} = \bar{y}_t$ where $\bar{y}_t = \bar{y}_{\tilde{R}}$ for all t under the fixed scheme and $\bar{y}_t = \tilde{R}^{-1} \sum_{s=t-\tilde{R}+1}^t y_s$ under the rolling scheme.

- Fixed scheme. For $t = 1, \dots, R$ set $y_t = 2\eta_t$ where $\eta_t = \varepsilon_t + \sum_{j=1}^{h-1} \theta_j \varepsilon_{t-j}$ with $\varepsilon_t \sim i.i.d.N(0, \sigma^2)$. For $t = R + 1, \dots, T$ set $y_t = \frac{1}{2}\bar{y}_R + \eta_t$.
- Rolling scheme. For $t = 1, \dots, T$ generate y_t as a stationary $ARMA(R + h - 1, h - 1)$ with moving average parameters θ_j , $j = 1, \dots, h - 1$, and autoregressive parameters α_j set so that the first $h - 1$ values are zero (i.e. $\alpha_j = 0$ for $j = 1, \dots, h - 1$) and the remaining R values are $\frac{1}{2R}$ (i.e. $\alpha_j = \frac{1}{2R}$ for $j = h, \dots, R + h - 1$).⁷
- As we did for the unconditional tests, we set $\sigma^2 = 1$, $\theta_j = (0.5)^j$, let $h = 1, 3, 12$, and use the same values of R , \tilde{R} , and P .
- We consider two test statistics. The first is simply the baseline t-statistic $P^{-1/2} \sum_{t=\tilde{R}}^{T-h} \hat{d}_{t+h} / \hat{\omega}$ which is asymptotically standard Normal under the null. But as noted in GW, since the null hypothesis implies \hat{d}_{t+h} is uncorrelated with any observable $z_t \in \mathfrak{S}_t$, $P^{-1/2} \sum_{t=\tilde{R}}^{T-h} \hat{d}_{t+h} z_t$ will also be asymptotically Normal under the null hypothesis. Since z_t is allowed to be non-scalar they suggest a Wald-statistic of the form $(P^{-1/2} \sum_{t=\tilde{R}}^{T-h} \hat{d}_{t+h} z_t)' \hat{\Omega}^{-1} (P^{-1/2} \sum_{t=\tilde{R}}^{T-h} \hat{d}_{t+h} z_t)$ which will be asymptotically

⁷The initial conditions are set to zero and a burn-in of 10,000 is used.

chi-square with $\dim(z_t)$ degrees of freedom under the null so long as $\hat{\Omega}$ is consistent for the long-run variance of $\hat{d}_{t+h}z_t$. Following the suggestion in GW, we investigate the usefulness of this statistic using $z_t = (1, \hat{d}_t)'$ as the instrument vector.

- When constructing $\hat{\Omega}$, we follow a slightly different approach than when we did for the unconditional tests. When $h = 1$, there is no serial correlation in either \hat{d}_{t+1} or $\hat{d}_{t+1}z_t$ and hence their variance is the long-run variance. We therefore construct $\hat{\Omega}$ as $P^{-1} \sum_{t=\tilde{R}}^{T-1} \hat{d}_{t+1}^2 z_t z_t'$ and $\hat{\omega}^2$ as $P^{-1} \sum_{t=\tilde{R}}^{T-1} \hat{d}_{t+1}^2$ when $h = 1$. For the longer horizons, we again used the Bartlett kernel and experimented with a bandwidth of h (as suggested by GW) and the data-dependent version. Since the actual size of the tests were weakly better using the data dependent rule we use that bandwidth when $h > 1$.

3.1 Fixed

In Table 3 we provide results associated with the scalar version of the test when the fixed scheme is used in the DGP and in the modeling. The layout of the table is the same as for the previous tables and hence rows in bold are actual size while those not in bold are actual power. In contrast to the results for the unconditional test, the fixed scheme can exhibit reasonable size in large enough samples. When $h = 1$, actual size is quite good for all sample sizes. But as we increase the forecast horizon, the need to estimate long-run variances again induces size distortions. Even so, it is comforting to see that for $h = 3$, the distortions diminish quickly as the sample size increases. At the longest horizon $h = 12$, the actual size improves as the sample increases but remains slightly elevated even when $P = 1000$.

Under the alternative in which $\tilde{R} \neq R$ but the fixed model is used to form forecasts (Table 3, not bold) power is poor with rejection frequencies generally well below 40% for all but the largest sample size. For the largest sample size, actual power can be as large as 75% but remains as low as 10% for some permutations of $\tilde{R} \neq R$. Holding (\tilde{R}, P) constant, in Table 4 we find that under the alternative in which the rolling model is used to form forecasts, actual power is typically the same as, but sometimes worse than in Table 3.

In Tables 5 and 6 we repeat the same experiments as in Tables 3 and 4 but now using the Wald form of the test with $z_t = (1, \hat{d}_t)'$ as the instrument vector. In Table 5, when $h = 1$, the actual size of the test is reasonable for all but the smallest out-of-sample size with rejection frequencies as high as 11% but more often ranging from 5% to 8%. At the longer horizons, actual size of the test is quite poor at the smaller sample sizes with

rejection frequencies as high as 45%. Even so, consistent with the theory, actual size improves monotonically as the sample size increases.

Under the alternative in which $\tilde{R} \neq R$, but the fixed model is used to form forecasts (Table 5, not bold), actual power is poor when $h = 1$ with rejection frequencies generally below 15% unless $\tilde{R} = 25$ or $P = 1000$ and $\tilde{R} < R$ in which case actual power can be as high as 70%. When $h > 1$, actual power is improved at all sample sizes particularly when $\tilde{R} < R$. Across all horizons there is a tendency for actual power to be higher for smaller values of \tilde{R} .

Under the alternative in which the rolling model is used to form forecasts (Table 6), actual power is a mixed bag. When $h = 1$, actual power is never above 20% unless $P = 1000$. When the sample size is large, the rejection frequencies are generally higher for smaller values of \tilde{R} and can get as high as 60%. Actual power is higher at the longer horizons holding (\tilde{R}, P) constant. As was the case for $h = 1$, actual power tends to rise as \tilde{R} becomes smaller with rejection frequencies approaching 70%. Somewhat oddly, there are times when the rejection frequencies are u-shaped as P increases for a fixed \tilde{R} . As an example, for $h = 3$, when $R = 25$ and $\tilde{R} = 75$ the rejection frequencies range from 28%, 21%, 20%, 20%, to 40% as P ranges from 25, 75, 125, 175, to 1000.

3.2 Rolling

In this section, the simulations parallel those in the previous subsection but when the rolling scheme is used to define the DGPs. In Table 7, we report rejection frequencies when both the DGP and model are based on the rolling scheme and the scalar test statistic is used for inference. Under the null, actual size of the test is quite good when $h = 1$. As the horizon increases, size distortions arise, especially when P is small but again, in alignment with the theory, actual size improves substantially as the sample size P increases with rejection frequencies ranging from 6% to 8% when $P = 1000$.

As it was for the fixed case in Table 3, under the alternative in which $\tilde{R} \neq R$ and the rolling model is used to form forecasts (Table 7, not bold) power is often poor with rejection frequencies typically below 20% for all but the largest sample size. For the largest sample size, actual power can be as large as 80% but remains as low as 5% for some permutations of $\tilde{R} \neq R$. Under the alternative in which the fixed model is used to form forecasts (Table 8), actual power is, for most permutations of (\tilde{R}, P) , substantially better than in Table 7.

In Tables 9 and 10 we repeat the same experiments as in Tables 7 and 8 but now using

the Wald form of the test with $z_t = (1, \hat{d}_t)'$ as the instrument vector. In Table 9, the actual size is modestly-to-severely oversized as it was in Table 5 for the fixed scheme. At the longer horizons, actual size of the test is typically poor at the smaller sample sizes with rejection frequencies as high as 40%. Even so, consistent with the theory, at all horizons actual size generally improves as the sample size increases.

Under the alternative in which the fixed model is used to form forecasts (Table 10), actual power again has a wide range. When $h = 1$, power approaches 40% for smaller permutations of R and \tilde{R} but are often below 20%. As the sample size P rises to 1000, actual power improves to as much as 70%. As the horizon increases, power increases almost uniformly holding (\tilde{R}, P) constant. For all horizons, actual power typically improves as \tilde{R} declines. And as was the case for the fixed DGP, there are odd instances in which the rejection frequencies are u-shaped for a fixed value of \tilde{R} as P increases.

In Tables 9 and 10 (as well as 5 and 6 under the fixed scheme) it is worth reiterating that the Bartlett kernel is being used to estimate the long-run variance of $\hat{d}_{t+h}z_t$. While this obviously captures the role of serial correlation in the long-run variance it also captures the presence of any conditional heteroskedasticity. This is particularly important given that \hat{d}_t is an element of z_t . To understand why, suppose $h = 1$, let $\tilde{R} = R$, and consider the second element of $P^{-1/2} \sum_{t=R}^{T-1} \hat{d}_{t+1}z_t$. Straightforward algebra reveals that $P^{-1/2} \sum_{t=R}^{T-1} \hat{d}_{t+1}\hat{d}_t$ is asymptotically Normal with zero mean and variance $\Omega_{22} = E\hat{d}_{t+1}^2\hat{d}_t^2$. Since $\hat{d}_{t+1} = 2\varepsilon_{t+1}\bar{y}_t$ and $\varepsilon_{t+1} \sim i.i.d.N(0, 1)$ we obtain $E\hat{d}_{t+1}^2\hat{d}_t^2 = 16E(\varepsilon_t^2\bar{y}_t^2\bar{y}_{t-1}^2) \neq 16E(\bar{y}_t^2)E(\bar{y}_{t-1}^2) = E\hat{d}_{t+1}^2E\hat{d}_t^2$ and thus \hat{d}_{t+1} exhibits conditional heteroskedasticity. In unreported results akin to those for $h = 1$ under the rolling scheme in Table 9, when we did not account for conditionally heteroskedasticity we obtained rejection frequencies ranging from 16% to 19% when $P = 1000$. But, as noted above, when we accounted for conditional heteroskedasticity the rejection frequencies were 5%, in line with the nominal size of the test.

4 Conclusion

In this paper we have modest goals. The first goal is simply to provide examples of DGPs that can be shown to satisfy each of the null hypotheses delineated in GW. With respect to the null hypothesis of equal unconditional, finite-sample predictive ability, Giacomini and Rossi (2010) provide a DGP that can satisfy the null but investigate size and power properties only in the context of the rolling scheme. There are no such simulations for the fixed scheme. In addition, there exist no simulations that delineate a DGP that satisfies

the null hypothesis of equal conditional, finite-sample predictive ability. In large part we suspect this arises because the null hypothesis of equal conditional finite-sample predictive ability is extremely narrow and imposes very strong restrictions on the DGP. As shown in GW, given two sequences of point forecasts $\hat{y}_{1,t+h}$ and $\hat{y}_{2,t+h}$ and assuming a quadratic loss function, the DGP for y must satisfy $y_{t+h} = \frac{1}{2}(\hat{y}_{1,t+h} + \hat{y}_{2,t+h}) + \eta_{t+h}$ for $\eta_{t+h} \sim MA(h-1)$ for this null hypothesis to hold. We have been able to delineate such a DGP but only for an application in which a no-change forecast is being compared to one formed using a location model. We have not been able to establish an example that permits a broader collection of regression models.

Nevertheless, given these DGPs, our second goal is to provide simulation-based evidence on the finite sample size and power of the tests suggested in GW. Our results only sometimes reinforce the conclusions provided in GW. The tests can be accurately sized given large enough samples, especially at shorter horizons. They too can exhibit substantial power when deviations from the null are large or the sample size is large. But by providing a DGP-based set of simulations we are able to diagnose details that are not particularly obvious in those provided by GW. The most important of which is simply that the proposed test statistic $P^{-1/2} \sum_{t=\tilde{R}}^{T-h} \hat{d}_{t+h} / \hat{\omega}$, is not asymptotically Normal when the fixed scheme is used and we are testing the null of equal unconditional finite-sample predictability. In fact, the test statistic diverges under the null since the loss differentials are not asymptotically independent. We also find that for this null, by construction, the rolling scheme induces very high orders of serial correlation in the loss differentials. As such, a large bandwidth is required when estimating the long-run variance which, unfortunately, introduces estimation error and causes size distortions. Finally, we find that when conducting the test of conditional predictive ability it is particularly important to account for conditional heteroskedasticity if a lagged loss differential is used as an instrument.

References

- Andrews, Donald W. K., and J. Christopher Monahan (1992), “An Improved Heteroskedasticity and Autocorrelation Consistent Covariance Matrix Estimator,” *Econometrica* 60, 953–966.
- 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 (2015), “Nested Forecast Model Comparisons: A New Approach to Testing Equal Accuracy,” *Journal of Econometrics* 186, 160–177.
- Coroneo, Laura, and Fabrizio Iacone (2018), “Comparing Predictive Accuracy in Small Samples Using Fixed-Smoothing Asymptotics,” manuscript.
- Corradi, Valentina, Olivetti, Claudia, and Norman R. Swanson (2001), “Predictive Ability With Cointegrated Variables,” *Journal of Econometrics* 104, 315–358.
- Diebold, Francis X. (2015), “Comparing Predictive Accuracy, Twenty Years Later: A Personal Perspective on the Use and Abuse of Diebold-Mariano Tests,” *Journal of Business and Economic Statistics* 33, 1–9.
- Diebold, Francis X., and Roberto S. Mariano (1995), “Comparing Predictive Accuracy,” *Journal of Business and Economic Statistics* 13, 253–263.
- Giacomini, Rafaella, and Barbara Rossi (2009), “Detecting and Predicting Forecast Breakdowns,” *Review of Economic Studies* 76, 669–705.
- Giacomini, Rafaella, and Barbara Rossi (2010), “Forecast Comparisons in Unstable Environments,” *Journal of Applied Econometrics* 25, 595–620.
- Giacomini, Rafaella, and Halbert White (2006), “Tests of Conditional Predictive Ability,” *Econometrica* 74, 1545–1578.
- Goncalves, Silvia, Perron, Benoit, and Michael W. McCracken (2017), “Tests of Equal Accuracy for Nested Models with Estimated Factors,” *Journal of Econometrics* 198, 231–252.
- Li, Jia, and Andrew Patton (2018), “Asymptotic Inference about Predictive Accuracy Using High Frequency Data,” *Journal of Econometrics* 203, 223–240.
- McCracken, Michael W. (2007), “Asymptotics for Out-of-Sample Tests of Granger Causality,” *Journal of Econometrics* 140, 719–752.
- Newey, Whitney K., and Kenneth D. West (1987). “A Simple, Positive Semi-definite, 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, 631–653.
- Timmermann, Allan and Yinchu Zhu (2017), “Tests of Forecasting Performance and Choice of Estimation Window,” manuscript.
- West, Kenneth D. (1996), “Asymptotic Inference About Predictive Ability,” *Econometrica*

64, 1067-1084.

White, Halbert (2000), "A Reality Check for Data Snooping," *Econometrica* 68, 1097-1127.

White, Halbert and Ian Domowitz (1984), "Nonlinear Regression with Dependent Observations," *Econometrica* 52, 142-162.

Tables

Table 1: Unconditional: Fixed Model

R	\tilde{R}	$h = 1$					$h = 3$					$h = 12$				
		P					P					P				
		25	75	125	175	1000	25	75	125	175	1000	25	75	125	175	1000
25	25	0.15	0.24	0.33	0.39	0.70	0.23	0.29	0.37	0.41	0.72	0.28	0.33	0.40	0.46	0.73
	75	0.13	0.19	0.26	0.32	0.74	0.20	0.26	0.31	0.37	0.74	0.25	0.28	0.34	0.40	0.75
	125	0.13	0.19	0.25	0.30	0.76	0.20	0.23	0.28	0.35	0.77	0.26	0.28	0.32	0.40	0.79
	175	0.12	0.17	0.24	0.30	0.78	0.20	0.22	0.29	0.34	0.80	0.26	0.27	0.32	0.39	0.80
75	25	0.13	0.18	0.25	0.28	0.62	0.20	0.23	0.29	0.34	0.63	0.24	0.28	0.32	0.36	0.66
	75	0.11	0.13	0.17	0.19	0.51	0.17	0.17	0.20	0.24	0.54	0.21	0.20	0.23	0.27	0.55
	125	0.10	0.11	0.14	0.17	0.48	0.17	0.15	0.19	0.22	0.52	0.21	0.20	0.21	0.25	0.54
	175	0.10	0.12	0.14	0.17	0.47	0.17	0.16	0.17	0.21	0.50	0.22	0.20	0.20	0.23	0.52
125	25	0.12	0.17	0.23	0.27	0.59	0.19	0.22	0.27	0.31	0.60	0.25	0.27	0.31	0.35	0.64
	75	0.11	0.11	0.14	0.18	0.44	0.18	0.16	0.19	0.21	0.48	0.22	0.20	0.21	0.23	0.50
	125	0.10	0.10	0.13	0.14	0.40	0.16	0.15	0.16	0.17	0.44	0.20	0.18	0.19	0.20	0.46
	175	0.10	0.10	0.12	0.13	0.38	0.17	0.14	0.15	0.16	0.40	0.21	0.18	0.18	0.20	0.44
175	25	0.12	0.17	0.22	0.27	0.57	0.20	0.22	0.26	0.31	0.59	0.25	0.25	0.28	0.33	0.62
	75	0.11	0.11	0.13	0.15	0.43	0.16	0.15	0.18	0.20	0.45	0.21	0.18	0.21	0.22	0.47
	125	0.10	0.10	0.12	0.13	0.35	0.17	0.14	0.15	0.17	0.41	0.21	0.18	0.17	0.20	0.42
	175	0.09	0.10	0.11	0.12	0.34	0.15	0.13	0.15	0.16	0.36	0.22	0.17	0.17	0.19	0.39

Notes: R is the window size used to generate the time series, \tilde{R} is the window size used to estimate the model parameters before generating a forecast, h is the forecast horizon, and P is the number of forecast periods. Each entry in the table represents the fraction of replications where the null hypothesis was rejected at the 5% level using the standard Normal critical values of a two-sided test for a given R , \tilde{R} , h , and P . Bolded rows indicate when the null hypothesis holds (i.e. when $R = \tilde{R}$). The test statistic takes the form shown in the second bullet of Section 2. The long-run variance used in the calculation of the test statistic was estimated using the Bartlett kernel with the bandwidth set to $\lfloor 4(P/100)^{2/9} \rfloor + 1$.

Table 2: Unconditional: Rolling Model

R	\tilde{R}	$h = 1$					$h = 3$					$h = 12$				
		P					P					P				
		25	75	125	175	1000	25	75	125	175	1000	25	75	125	175	1000
25	25	0.07	0.05	0.05	0.04	0.03	0.12	0.07	0.06	0.05	0.04	0.20	0.11	0.10	0.09	0.07
	75	0.08	0.07	0.06	0.05	0.41	0.16	0.11	0.09	0.10	0.45	0.22	0.16	0.14	0.15	0.48
	125	0.10	0.09	0.10	0.12	0.67	0.17	0.14	0.15	0.16	0.71	0.24	0.20	0.20	0.20	0.73
	175	0.10	0.12	0.14	0.15	0.79	0.18	0.16	0.17	0.20	0.81	0.24	0.22	0.23	0.24	0.82
75	25	0.09	0.11	0.12	0.15	0.51	0.12	0.12	0.13	0.17	0.50	0.21	0.16	0.16	0.18	0.47
	75	0.08	0.06	0.06	0.05	0.03	0.13	0.09	0.07	0.06	0.03	0.19	0.12	0.10	0.08	0.06
	125	0.08	0.06	0.05	0.04	0.04	0.13	0.09	0.08	0.06	0.05	0.19	0.13	0.10	0.09	0.07
	175	0.09	0.06	0.06	0.05	0.06	0.15	0.10	0.08	0.08	0.08	0.20	0.13	0.12	0.10	0.11
125	25	0.10	0.13	0.17	0.20	0.72	0.12	0.14	0.17	0.20	0.70	0.19	0.18	0.18	0.23	0.67
	75	0.09	0.08	0.06	0.07	0.08	0.13	0.09	0.08	0.08	0.10	0.19	0.12	0.11	0.10	0.12
	125	0.08	0.07	0.05	0.05	0.03	0.12	0.09	0.07	0.07	0.04	0.18	0.13	0.10	0.08	0.05
	175	0.09	0.07	0.06	0.05	0.02	0.13	0.10	0.08	0.07	0.03	0.19	0.13	0.10	0.11	0.05
175	25	0.10	0.14	0.20	0.23	0.82	0.13	0.14	0.20	0.23	0.80	0.20	0.18	0.22	0.26	0.74
	75	0.09	0.08	0.08	0.08	0.15	0.13	0.10	0.09	0.10	0.16	0.19	0.13	0.10	0.11	0.19
	125	0.08	0.07	0.07	0.06	0.06	0.14	0.09	0.08	0.07	0.07	0.19	0.12	0.10	0.09	0.08
	175	0.09	0.07	0.06	0.06	0.03	0.13	0.10	0.08	0.08	0.04	0.19	0.12	0.11	0.09	0.05

Notes: See notes to Table 1.

Table 3: Conditional: Fixed DGP and Fixed Model, Scalar Form

R	\tilde{R}	$h = 1$					$h = 3$					$h = 12$				
		P					P					P				
		25	75	125	175	1000	25	75	125	175	1000	25	75	125	175	1000
25	25	0.06	0.06	0.06	0.05	0.05	0.14	0.10	0.08	0.08	0.07	0.19	0.13	0.11	0.10	0.08
	75	0.09	0.11	0.15	0.18	0.48	0.16	0.17	0.20	0.22	0.51	0.23	0.21	0.22	0.25	0.52
	125	0.08	0.12	0.17	0.21	0.51	0.17	0.19	0.23	0.25	0.54	0.23	0.22	0.26	0.28	0.55
	175	0.09	0.13	0.18	0.22	0.52	0.18	0.19	0.22	0.26	0.55	0.22	0.23	0.25	0.29	0.57
75	25	0.10	0.18	0.28	0.37	0.70	0.18	0.25	0.35	0.41	0.73	0.23	0.32	0.39	0.46	0.73
	75	0.06	0.05	0.05	0.05	0.05	0.15	0.10	0.08	0.08	0.06	0.19	0.13	0.11	0.10	0.08
	125	0.06	0.07	0.07	0.07	0.20	0.15	0.12	0.12	0.11	0.23	0.19	0.14	0.14	0.13	0.25
	175	0.07	0.07	0.08	0.09	0.25	0.16	0.12	0.12	0.12	0.27	0.20	0.16	0.14	0.15	0.32
125	25	0.09	0.14	0.22	0.31	0.71	0.18	0.20	0.27	0.36	0.73	0.25	0.24	0.33	0.40	0.75
	75	0.07	0.10	0.15	0.19	0.45	0.16	0.17	0.21	0.23	0.48	0.21	0.22	0.25	0.27	0.50
	125	0.07	0.06	0.05	0.05	0.05	0.15	0.11	0.09	0.09	0.07	0.19	0.14	0.11	0.10	0.09
	175	0.07	0.06	0.06	0.06	0.13	0.15	0.10	0.09	0.10	0.14	0.20	0.13	0.12	0.13	0.17
175	25	0.09	0.13	0.20	0.26	0.70	0.18	0.20	0.25	0.32	0.71	0.24	0.24	0.27	0.35	0.73
	75	0.07	0.09	0.11	0.16	0.51	0.17	0.14	0.17	0.22	0.53	0.21	0.17	0.20	0.25	0.55
	125	0.07	0.08	0.11	0.14	0.32	0.16	0.14	0.15	0.19	0.34	0.19	0.18	0.18	0.20	0.37
	175	0.06	0.05	0.05	0.06	0.05	0.14	0.09	0.09	0.08	0.06	0.19	0.13	0.11	0.11	0.07

Notes: See notes to Table 1. The long-run variance used in the calculation of the test statistic was estimated using the Bartlett kernel with the bandwidth set to 1 in the case of $h = 1$ and $\lfloor 4(P/100)^{2/9} \rfloor + 1$ otherwise.

Table 4: Conditional: Fixed DGP and Rolling Model, Scalar Form

R	\tilde{R}	$h = 1$					$h = 3$					$h = 12$				
		P					P					P				
		25	75	125	175	1000	25	75	125	175	1000	25	75	125	175	1000
25	25	0.06	0.08	0.11	0.14	0.54	0.13	0.13	0.16	0.20	0.57	0.20	0.17	0.19	0.23	0.56
	75	0.07	0.08	0.11	0.14	0.42	0.15	0.13	0.17	0.19	0.44	0.21	0.18	0.18	0.22	0.46
	125	0.07	0.11	0.13	0.16	0.42	0.16	0.16	0.18	0.20	0.44	0.22	0.21	0.21	0.23	0.46
	175	0.07	0.11	0.15	0.18	0.44	0.17	0.18	0.19	0.22	0.47	0.22	0.21	0.24	0.25	0.48
75	25	0.07	0.12	0.15	0.19	0.63	0.13	0.18	0.20	0.25	0.67	0.20	0.21	0.22	0.27	0.66
	75	0.06	0.05	0.06	0.06	0.25	0.13	0.10	0.10	0.10	0.30	0.19	0.13	0.12	0.12	0.32
	125	0.05	0.06	0.06	0.07	0.22	0.13	0.10	0.10	0.10	0.25	0.19	0.14	0.12	0.12	0.27
	175	0.06	0.07	0.06	0.07	0.22	0.14	0.10	0.10	0.11	0.25	0.19	0.15	0.13	0.13	0.26
125	25	0.07	0.12	0.17	0.23	0.69	0.14	0.17	0.24	0.28	0.71	0.20	0.20	0.26	0.30	0.70
	75	0.06	0.09	0.09	0.10	0.29	0.13	0.11	0.13	0.14	0.32	0.18	0.15	0.16	0.16	0.35
	125	0.06	0.06	0.05	0.06	0.17	0.13	0.09	0.09	0.09	0.20	0.17	0.12	0.11	0.11	0.21
	175	0.06	0.05	0.06	0.05	0.14	0.13	0.09	0.09	0.09	0.17	0.18	0.13	0.11	0.11	0.20
175	25	0.07	0.12	0.18	0.25	0.72	0.14	0.18	0.25	0.32	0.75	0.20	0.20	0.25	0.33	0.74
	75	0.06	0.07	0.09	0.11	0.30	0.13	0.11	0.13	0.14	0.35	0.18	0.14	0.15	0.18	0.36
	125	0.06	0.07	0.08	0.08	0.20	0.13	0.11	0.10	0.11	0.24	0.18	0.13	0.12	0.14	0.25
	175	0.06	0.05	0.05	0.05	0.13	0.13	0.08	0.09	0.08	0.15	0.18	0.12	0.12	0.11	0.18

Notes: See notes to Table 1. The long-run variance used in the calculation of the test statistic was estimated using the Bartlett kernel with the bandwidth set to 1 in the case of $h = 1$ and $\lfloor 4(P/100)^{2/9} \rfloor + 1$ otherwise. No values are bolded because the null hypothesis does not hold for any entry.

Table 5: Conditional: Fixed DGP and Fixed Model, Wald Form

R	\tilde{R}	$h = 1$					$h = 3$					$h = 12$				
		P					P					P				
		25	75	125	175	1000	25	75	125	175	1000	25	75	125	175	1000
25	25	0.11	0.06	0.07	0.06	0.05	0.28	0.14	0.12	0.11	0.06	0.42	0.19	0.15	0.13	0.09
	75	0.11	0.11	0.13	0.15	0.43	0.29	0.21	0.21	0.23	0.46	0.43	0.24	0.23	0.25	0.47
	125	0.13	0.13	0.14	0.18	0.47	0.30	0.20	0.23	0.24	0.48	0.44	0.25	0.25	0.27	0.51
	175	0.13	0.13	0.15	0.19	0.48	0.29	0.22	0.23	0.25	0.51	0.45	0.26	0.27	0.28	0.51
75	25	0.13	0.17	0.26	0.33	0.69	0.30	0.28	0.36	0.41	0.69	0.45	0.33	0.40	0.44	0.71
	75	0.11	0.08	0.05	0.06	0.06	0.26	0.15	0.12	0.11	0.07	0.44	0.18	0.15	0.12	0.08
	125	0.10	0.08	0.07	0.08	0.17	0.27	0.16	0.15	0.13	0.19	0.43	0.20	0.17	0.16	0.22
	175	0.11	0.08	0.08	0.08	0.21	0.28	0.16	0.15	0.14	0.24	0.44	0.21	0.17	0.16	0.27
125	25	0.13	0.13	0.19	0.27	0.68	0.31	0.22	0.29	0.36	0.71	0.44	0.26	0.32	0.39	0.73
	75	0.12	0.11	0.14	0.16	0.40	0.28	0.21	0.22	0.23	0.46	0.45	0.25	0.25	0.26	0.46
	125	0.11	0.07	0.06	0.05	0.05	0.27	0.15	0.11	0.10	0.08	0.43	0.19	0.14	0.13	0.09
	175	0.10	0.08	0.07	0.06	0.11	0.27	0.16	0.14	0.11	0.14	0.44	0.19	0.15	0.13	0.15
175	25	0.13	0.13	0.16	0.23	0.66	0.30	0.23	0.25	0.30	0.69	0.45	0.27	0.27	0.33	0.71
	75	0.11	0.09	0.10	0.16	0.47	0.28	0.18	0.19	0.21	0.49	0.45	0.22	0.22	0.25	0.52
	125	0.12	0.10	0.11	0.13	0.27	0.28	0.19	0.19	0.19	0.31	0.44	0.22	0.20	0.22	0.33
	175	0.11	0.06	0.06	0.06	0.05	0.27	0.14	0.12	0.11	0.07	0.45	0.18	0.15	0.14	0.08

Notes: See notes to Table 1. The test statistic used takes the form of the Wald-statistic described in the fourth bullet of Section 3. Consequently, $\chi^2(2)$ critical values are used for inference (as opposed to the standard Normal critical values used in Table 1). The long-run variance used in the calculation of the test statistic was estimated using the Bartlett kernel with the bandwidth set to 1 in the case of $h = 1$ and $\lfloor 4((P - h)/100)^{2/9} \rfloor + 1$ otherwise. Note that h observations were lost because the instrument vector in the test statistic includes an h -lagged loss differential.

Table 6: Conditional: Fixed DGP and Rolling Model, Wald Form

R	\tilde{R}	$h = 1$					$h = 3$					$h = 12$				
		P					P					P				
		25	75	125	175	1000	25	75	125	175	1000	25	75	125	175	1000
25	25	0.12	0.10	0.11	0.13	0.48	0.25	0.18	0.19	0.22	0.55	0.40	0.21	0.20	0.23	0.53
	75	0.12	0.10	0.11	0.13	0.36	0.28	0.21	0.20	0.20	0.40	0.42	0.23	0.22	0.21	0.43
	125	0.13	0.11	0.13	0.15	0.38	0.27	0.20	0.20	0.22	0.43	0.41	0.24	0.24	0.23	0.43
	175	0.12	0.12	0.14	0.16	0.39	0.28	0.21	0.22	0.22	0.43	0.44	0.25	0.25	0.25	0.45
75	25	0.11	0.13	0.15	0.16	0.57	0.23	0.21	0.23	0.24	0.62	0.38	0.22	0.21	0.22	0.59
	75	0.11	0.08	0.08	0.07	0.20	0.25	0.16	0.15	0.15	0.26	0.40	0.18	0.15	0.15	0.27
	125	0.12	0.09	0.08	0.07	0.18	0.27	0.15	0.15	0.13	0.22	0.40	0.18	0.16	0.15	0.24
	175	0.11	0.08	0.08	0.08	0.18	0.26	0.16	0.15	0.14	0.23	0.41	0.18	0.17	0.16	0.25
125	25	0.12	0.12	0.15	0.19	0.59	0.23	0.20	0.25	0.27	0.67	0.37	0.20	0.23	0.25	0.64
	75	0.11	0.09	0.11	0.10	0.24	0.24	0.17	0.18	0.17	0.29	0.39	0.18	0.17	0.15	0.30
	125	0.11	0.07	0.07	0.07	0.13	0.25	0.15	0.15	0.13	0.19	0.39	0.17	0.14	0.14	0.19
	175	0.11	0.07	0.07	0.07	0.11	0.26	0.15	0.14	0.12	0.15	0.39	0.18	0.15	0.14	0.19
175	25	0.11	0.12	0.15	0.20	0.62	0.23	0.20	0.26	0.29	0.69	0.38	0.20	0.22	0.27	0.66
	75	0.12	0.08	0.09	0.09	0.23	0.24	0.17	0.16	0.17	0.31	0.39	0.18	0.16	0.17	0.29
	125	0.12	0.09	0.08	0.09	0.17	0.25	0.16	0.16	0.16	0.21	0.40	0.17	0.17	0.15	0.21
	175	0.11	0.07	0.07	0.06	0.11	0.25	0.15	0.14	0.13	0.14	0.40	0.19	0.14	0.13	0.16

Notes: See notes to Table 5. No values are bolded because the null hypothesis does not hold for any entry.

Table 7: Conditional: Rolling DGP and Rolling Model, Scalar Form

R	\tilde{R}	$h = 1$					$h = 3$					$h = 12$				
		P					P					P				
		25	75	125	175	1000	25	75	125	175	1000	25	75	125	175	1000
25	25	0.05	0.06	0.05	0.05	0.05	0.12	0.09	0.08	0.07	0.07	0.17	0.10	0.09	0.08	0.07
	75	0.11	0.10	0.08	0.08	0.07	0.18	0.14	0.11	0.10	0.07	0.20	0.15	0.12	0.11	0.09
	125	0.13	0.14	0.13	0.12	0.10	0.22	0.20	0.16	0.15	0.10	0.23	0.20	0.17	0.16	0.12
	175	0.14	0.18	0.16	0.14	0.12	0.24	0.23	0.21	0.18	0.12	0.25	0.24	0.20	0.19	0.15
75	25	0.07	0.09	0.10	0.15	0.49	0.12	0.14	0.16	0.20	0.53	0.18	0.15	0.18	0.20	0.50
	75	0.06	0.05	0.05	0.04	0.05	0.13	0.10	0.08	0.08	0.06	0.17	0.11	0.09	0.08	0.08
	125	0.07	0.07	0.06	0.07	0.05	0.15	0.11	0.09	0.10	0.06	0.18	0.14	0.11	0.09	0.08
	175	0.08	0.08	0.09	0.08	0.07	0.17	0.14	0.13	0.12	0.08	0.20	0.16	0.15	0.13	0.08
125	25	0.07	0.10	0.13	0.17	0.68	0.13	0.15	0.20	0.24	0.71	0.19	0.18	0.20	0.23	0.66
	75	0.06	0.06	0.07	0.08	0.14	0.13	0.10	0.10	0.10	0.16	0.17	0.12	0.11	0.12	0.17
	125	0.06	0.05	0.05	0.05	0.05	0.14	0.10	0.09	0.08	0.06	0.19	0.11	0.09	0.10	0.07
	175	0.07	0.07	0.07	0.05	0.05	0.15	0.11	0.10	0.09	0.07	0.19	0.13	0.12	0.11	0.08
175	25	0.07	0.10	0.14	0.20	0.77	0.13	0.15	0.19	0.25	0.80	0.19	0.18	0.21	0.26	0.74
	75	0.06	0.06	0.07	0.08	0.19	0.14	0.10	0.10	0.11	0.23	0.18	0.12	0.12	0.13	0.25
	125	0.07	0.06	0.05	0.06	0.09	0.14	0.10	0.09	0.10	0.11	0.19	0.13	0.11	0.10	0.12
	175	0.07	0.06	0.05	0.06	0.05	0.14	0.10	0.08	0.07	0.06	0.19	0.13	0.10	0.08	0.07

Notes: See notes to Table 1. The long-run variance used in the calculation of the test statistic was estimated using the Bartlett kernel with the bandwidth set to 1 in the case of $h = 1$ and $\lfloor 4(P/100)^{2/9} \rfloor + 1$ otherwise.

Table 8: Conditional: Rolling DGP and Fixed Model, Scalar Form

R	\tilde{R}	$h = 1$					$h = 3$					$h = 12$				
		P					P					P				
		25	75	125	175	1000	25	75	125	175	1000	25	75	125	175	1000
25	25	0.16	0.29	0.36	0.43	0.67	0.27	0.36	0.41	0.45	0.66	0.26	0.36	0.40	0.44	0.65
	75	0.16	0.30	0.35	0.38	0.57	0.28	0.36	0.37	0.39	0.58	0.26	0.35	0.39	0.43	0.60
	125	0.17	0.28	0.33	0.35	0.55	0.28	0.35	0.36	0.38	0.52	0.28	0.35	0.38	0.40	0.56
	175	0.18	0.28	0.33	0.33	0.51	0.28	0.35	0.37	0.38	0.49	0.27	0.34	0.38	0.39	0.53
75	25	0.10	0.17	0.27	0.33	0.63	0.20	0.25	0.31	0.37	0.62	0.23	0.27	0.33	0.38	0.64
	75	0.08	0.15	0.21	0.27	0.52	0.17	0.21	0.27	0.31	0.54	0.20	0.22	0.27	0.31	0.55
	125	0.09	0.15	0.21	0.26	0.49	0.18	0.22	0.26	0.30	0.50	0.21	0.24	0.29	0.32	0.52
	175	0.10	0.16	0.23	0.27	0.47	0.18	0.22	0.27	0.31	0.48	0.21	0.25	0.29	0.32	0.50
125	25	0.10	0.16	0.21	0.28	0.60	0.20	0.22	0.27	0.33	0.60	0.23	0.25	0.30	0.36	0.62
	75	0.08	0.11	0.15	0.20	0.49	0.17	0.16	0.20	0.24	0.49	0.20	0.19	0.22	0.26	0.51
	125	0.07	0.10	0.14	0.18	0.45	0.16	0.16	0.20	0.23	0.46	0.20	0.20	0.22	0.26	0.48
	175	0.07	0.11	0.15	0.19	0.43	0.17	0.17	0.19	0.25	0.44	0.19	0.19	0.22	0.26	0.47
175	25	0.10	0.16	0.22	0.26	0.60	0.19	0.22	0.27	0.30	0.61	0.24	0.25	0.28	0.32	0.62
	75	0.08	0.10	0.12	0.15	0.46	0.17	0.15	0.16	0.21	0.49	0.21	0.18	0.20	0.23	0.49
	125	0.07	0.07	0.12	0.14	0.41	0.16	0.13	0.16	0.18	0.44	0.21	0.16	0.18	0.20	0.46
	175	0.07	0.08	0.11	0.15	0.40	0.16	0.14	0.15	0.20	0.42	0.20	0.16	0.18	0.21	0.43

Notes: See notes to Table 1. The long-run variance used in the calculation of the test statistic was estimated using the Bartlett kernel with the bandwidth set to 1 in the case of $h = 1$ and $\lfloor 4(P/100)^{2/9} \rfloor + 1$ otherwise. No values are bolded because the null hypothesis does not hold for any entry.

Table 9: Conditional: Rolling DGP and Rolling Model, Wald Form

R	\tilde{R}	$h = 1$					$h = 3$					$h = 12$				
		P					P					P				
		25	75	125	175	1000	25	75	125	175	1000	25	75	125	175	1000
25	25	0.11	0.08	0.08	0.07	0.05	0.24	0.15	0.15	0.13	0.09	0.37	0.16	0.12	0.11	0.09
	75	0.15	0.11	0.09	0.09	0.08	0.30	0.20	0.17	0.15	0.11	0.41	0.21	0.17	0.15	0.11
	125	0.17	0.14	0.13	0.12	0.12	0.32	0.24	0.21	0.18	0.16	0.43	0.25	0.20	0.18	0.18
	175	0.17	0.17	0.16	0.15	0.14	0.33	0.26	0.24	0.22	0.19	0.42	0.27	0.25	0.22	0.21
75	25	0.12	0.10	0.11	0.13	0.40	0.23	0.20	0.20	0.21	0.46	0.38	0.17	0.17	0.19	0.41
	75	0.11	0.08	0.07	0.06	0.05	0.25	0.16	0.13	0.13	0.08	0.40	0.16	0.13	0.12	0.08
	125	0.11	0.09	0.08	0.07	0.06	0.26	0.17	0.15	0.13	0.09	0.41	0.18	0.15	0.13	0.09
	175	0.12	0.11	0.09	0.09	0.06	0.28	0.19	0.17	0.15	0.09	0.41	0.20	0.17	0.15	0.10
125	25	0.12	0.10	0.12	0.14	0.58	0.24	0.19	0.21	0.23	0.63	0.38	0.19	0.18	0.21	0.57
	75	0.11	0.08	0.08	0.08	0.11	0.26	0.17	0.15	0.14	0.16	0.40	0.18	0.15	0.14	0.16
	125	0.11	0.07	0.07	0.06	0.05	0.26	0.16	0.14	0.12	0.08	0.40	0.17	0.14	0.12	0.09
	175	0.12	0.08	0.08	0.07	0.05	0.27	0.16	0.14	0.13	0.09	0.40	0.19	0.15	0.13	0.09
175	25	0.12	0.11	0.13	0.16	0.67	0.23	0.20	0.21	0.23	0.72	0.38	0.19	0.20	0.22	0.66
	75	0.11	0.08	0.08	0.08	0.14	0.26	0.15	0.15	0.14	0.20	0.40	0.18	0.15	0.14	0.21
	125	0.11	0.08	0.08	0.07	0.08	0.25	0.17	0.14	0.14	0.12	0.41	0.18	0.14	0.13	0.12
	175	0.11	0.08	0.07	0.07	0.05	0.25	0.17	0.13	0.12	0.09	0.41	0.18	0.14	0.13	0.08

Notes: See notes to Table 5.

Table 10: Conditional: Rolling DGP and Fixed Model, Wald Form

R	\tilde{R}	$h = 1$					$h = 3$					$h = 12$				
		P					P					P				
		25	75	125	175	1000	25	75	125	175	1000	25	75	125	175	1000
25	25	0.17	0.27	0.34	0.39	0.70	0.35	0.37	0.43	0.45	0.75	0.43	0.39	0.41	0.45	0.73
	75	0.19	0.26	0.32	0.36	0.63	0.35	0.37	0.40	0.43	0.70	0.45	0.36	0.38	0.43	0.68
	125	0.18	0.25	0.30	0.32	0.58	0.36	0.37	0.37	0.41	0.66	0.46	0.36	0.38	0.40	0.66
	175	0.19	0.25	0.30	0.31	0.54	0.36	0.36	0.37	0.39	0.63	0.45	0.36	0.37	0.40	0.63
75	25	0.14	0.16	0.23	0.29	0.60	0.30	0.26	0.31	0.35	0.62	0.43	0.29	0.32	0.37	0.63
	75	0.12	0.14	0.20	0.23	0.48	0.29	0.24	0.27	0.30	0.51	0.44	0.27	0.29	0.32	0.54
	125	0.12	0.14	0.19	0.23	0.45	0.29	0.24	0.27	0.29	0.50	0.44	0.26	0.29	0.32	0.50
	175	0.12	0.15	0.19	0.24	0.42	0.29	0.24	0.27	0.30	0.48	0.44	0.25	0.30	0.31	0.49
125	25	0.13	0.15	0.20	0.25	0.55	0.30	0.24	0.26	0.31	0.59	0.45	0.26	0.28	0.33	0.61
	75	0.11	0.11	0.14	0.17	0.45	0.28	0.19	0.21	0.24	0.46	0.42	0.23	0.23	0.27	0.48
	125	0.11	0.10	0.12	0.16	0.39	0.28	0.20	0.22	0.24	0.44	0.42	0.22	0.23	0.25	0.46
	175	0.12	0.10	0.13	0.16	0.38	0.29	0.20	0.20	0.24	0.42	0.44	0.23	0.23	0.25	0.44
175	25	0.13	0.14	0.18	0.22	0.56	0.29	0.23	0.25	0.29	0.57	0.45	0.27	0.27	0.32	0.58
	75	0.11	0.10	0.11	0.14	0.40	0.28	0.20	0.18	0.21	0.45	0.45	0.22	0.20	0.23	0.45
	125	0.12	0.09	0.11	0.13	0.36	0.27	0.18	0.17	0.19	0.40	0.42	0.21	0.20	0.21	0.41
	175	0.12	0.09	0.11	0.13	0.34	0.28	0.17	0.17	0.19	0.39	0.44	0.22	0.20	0.20	0.41

Notes: See notes to Table 5. No values are bolded because the null hypothesis does not hold for any entry.