

Alternative Tests for Correct Specification of Conditional Predictive Densities

Barbara Rossi¹ and Tatevik Sekhposyan²

July 14, 2017

Abstract

We propose a new framework for evaluating predictive densities in an environment where the estimation error of the parameters used to construct the densities is preserved asymptotically under the null hypothesis. The tests offer a simple way to evaluate the correct specification of predictive densities, where both the model specification and its estimation technique are evaluated jointly. Monte Carlo simulation results indicate that our tests are well sized and have good power in detecting misspecification. An empirical application to density forecasts of the Survey of Professional Forecasters shows the usefulness of our methodology.³

Keywords: Predictive Density, Probability Integral Transform, Kolmogorov-Smirnov Test, Cramér-von Mises Test, Forecast Evaluation

J.E.L. Codes: C22, C52, C53

¹ICREA-Universitat Pompeu Fabra, Barcelona GSE and CREI, c/Ramon Trias Fargas 25/27, Barcelona 08005, Spain; tel.: +34-93-542-1655; e-mail: barbara.rossi@upf.edu

²Texas A&M University, 3060 Allen Building, 4228 TAMU, College Station, TX 77843, USA; tel.: +1-979-862-8857; e-mail: tsekposyan@tamu.edu

³Acknowledgments: We thank T. Clark, F. Diebold, G. Ganics, A. Inoue, A. Patton, B. Perron, F. Ravazzolo, N. Swanson, M. Watson, seminar participants at George Washington Univ., Lehigh Univ., Univ. of California-Riverside, University of Mississippi, Texas A&M, Rutgers, CORE Louvain-la-Neuve, Univ. of New South Wales, Monash, ECB and the 2012 Econometric Society Australasian Meetings, the 2012 Time Series Econometrics Workshop in Zaragoza, the 2013 CIREQ Time Series and Financial Econometrics Conference, the 2013 St. Louis Fed Applied Time Series Econometrics Workshop, the 2013 UCL Conference on Frontiers in Macroeconometrics, the 2013 Conference on Forecasting Structure and Time Varying Parameter Patterns, the 2014 EC² Conference, the 2015 CAMP Workshop on Empirical Macroeconomics and the 2015 CIRANO-CIREQ Workshop on Data Revision in Macroeconomic Forecasting and Policy for comments. B. Rossi gratefully acknowledges financial support from the European Research Agency's Marie Curie Grant 303434, ERC Grant 615608 and the Spanish Ministry of Economy and Competitiveness, Grant ECO2015-68136-P, FEDER, UE.

1 Introduction

Policy institutions are becoming interested in complementing point forecasts with an accurate description of uncertainty. For instance, they are interested not only in knowing whether inflation is below its target, but also in understanding whether the realized inflation rate was forecasted to be a low probability event *ex-ante*. In fact, if researchers underestimate the uncertainty around point forecasts, it is possible that an event with a fairly high likelihood of occurrence is forecasted to be a very low probability event. An accurate description of uncertainty is therefore important in the decision making process of economic agents and policymakers. The interest in density forecasting has emerged in the survey by Elliott and Timmermann (2008) as well as in their recent book (Elliott and Timmermann, 2016), and has inspired several empirical contributions that have proposed new approaches to improve the forecasting performance of predictive densities, e.g. Aastveit, Foroni and Ravazzolo (2017), Ravazzolo and Vahey (2014) and Billio, Casarin, Ravazzolo and van Dijk (2013). The objective of this paper is to provide reliable tools for evaluating whether the uncertainty around point forecasts, and predictive densities in general, are correctly specified.

Many central banks periodically report fan charts to evaluate and communicate the uncertainty around point forecasts (e.g., see the various issues of the Bank of England Inflation Report or the Economic Bulletin of the Bank of Italy).⁴ Fan charts provide percentiles of the forecast distribution for macroeconomic variables of interest. Typically, central banks' fan charts are the result of convoluted methodologies that involve a variety of models and subjective assessments, although fan charts can be based on specific models as well.⁵ Figure 1 plots fan charts for US output growth (left panel) and the Federal Funds rate (right panel) based on a representative macroeconomic model by Smets and Wouters (2007) widely used in academia and policymaking.⁶ The fan charts display model-based forecasts made in 2000:IV for the next four quarters. The shaded areas in the figures depict the deciles of the forecast distribution and provide a visual impression of the uncertainty around the point forecasts (in this case, the median, marked by a solid line). Over the four quarterly horizons, uncertainty about output growth and interest rate forecasts has a very different pattern: the uncertainty surrounding output growth forecasts is constant across horizons,

⁴These publications are available at <http://www.bankofengland.co.uk/publications/Pages/inflationreport> and <https://www.bancaditalia.it/pubblicazioni/econo/bollec>, respectively.

⁵See for instance Clements (2004) for a discussion on the Bank of England fan charts.

⁶For a discussion on the forecasting ability of DSGE models, see Edge and Gürkaynak (2010), Edge, Kiley and Laforte (2010) and Gürkaynak et al. (2013).

while it depends on the horizon for interest rates. The dark, dash-dotted line in the figures plots the actual realized value of the target variable. Clearly, forecasts of interest rates are very imprecise (the realization is outside every forecast decile except for one-quarter-ahead horizon), whereas the model predicts output growth more accurately. Evaluating model-based forecast distributions amounts to understanding whether the model’s description of uncertainty was inaccurate or whether the realized values were indeed low probability events.

INSERT FIGURE 1 HERE

The methodologies that are currently available test whether the empirical distribution belongs to a given parametric density family with parameters evaluated at their pseudo-true values. Our paper derives new tools to evaluate whether predictive densities are correctly specified by focusing on evaluating their actual forecasting ability at models’ estimated parameter values, which, we argue, is more empirically useful to measure models’ actual forecasting ability in finite samples. In other words, we test whether the predictive densities are correctly specified given the parametric model and the estimation technique specified by the researcher. Accordingly, our tests do not require an asymptotic correction for parameter estimation error. Furthermore, our null hypothesis is that of correct specification of the density forecast, which, as we clarify in an example, can still hold even if the forecasting model is dynamically misspecified. Thus, even in the presence of dynamic misspecification, we obtain limiting distributions that are nuisance parameter free for one-step-ahead density forecasts. However, we also extend our framework to multiple-step-ahead forecasts, where the asymptotic distribution of our proposed tests is not nuisance parameter free.

Our approach, where parameter estimation error is maintained under the null hypothesis, is inspired by Amisano and Giacomini (2007). However, our approach is very different, as the latter focus on model selection by comparing the relative performance of competing models’ predictive densities, whereas we focus on evaluating the absolute performance of a model’s predictive density. Maintaining parameter estimation error under the null hypothesis has two advantages: (i) there is no need to correct the asymptotic distribution of the test statistics for parameter estimation error; and (ii) the asymptotic distribution of the test statistics is nuisance parameter free even when the model is dynamically misspecified.⁷ We derive our tests within the class of Kolmogorov-Smirnov and Cramér-von Mises-type tests commonly

⁷Note that (i) is not unique to cases where parameter estimation error is maintained under the null hypothesis; in fact, it also holds when parameter estimation error is asymptotically irrelevant, or when one uses martingalization techniques, as in Bai (2003).

used in the literature and show that our proposed tests have good size properties in small samples.

There are several existing approaches for testing the correct specification of a parametric density in-sample (e.g. Bai, 2003, Hong and Li, 2005, Corradi and Swanson, 2006a).⁸ Our paper focuses instead on the out-of-sample evaluation of predictive densities. The difference between in-sample and out-of-sample evaluation is that a model may fit well in-sample, and yet its out-of-sample forecasts may be poor (for example, if the distribution of the error changes between the in-sample estimation period and the out-of-sample evaluation period, or if the researcher overfitted the relevant distributional parameters). As such, our paper is related to a series of contributions which test whether observed realizations could have been generated by a given predictive distribution. Diebold et al. (1998, 1999) introduced the probability integral transform (PIT) into economics as a tool to test whether the empirical predictive distribution of surveys or empirical models matches the true, unobserved distribution that generates the data. Their approach tests for properties of the PITs, such as independence and uniformity, by treating the forecasts as primitive data, that is without correcting for estimation uncertainty associated with those forecasts.

Additional approaches proposed in the literature for assessing the correct calibration of predictive densities are the raw-moment-based test by Knueppel (2015), the likelihood ratio test by Berkowitz (2001), the non-parametric approach by Hong, Li and Zhao (2007), the bootstrap introduced by Corradi and Swanson (2006b,c) and the graphical devices by González-Rivera and Sun (2014).⁹ The null hypothesis in Hong, Li and Zhao (2007) and Corradi and Swanson (2006b,c) is that of correct specification of the density forecast at the pseudo-true (limiting) parameter values. Although this framework enables predictive density evaluation when the models are dynamically misspecified, it does not necessarily capture the actual measure of predictive ability that researchers are interested in, as in small samples the pseudo-true parameter values may not be representative of the actual predictive ability of the regressors. In the approach we propose, the main test statistic is the same as Corradi and Swanson's (2006a) one, although the null hypothesis is very different: it evaluates density forecasts at the estimated parameter values (as opposed to their population values). Thus,

⁸See also Bai and Ng (2005) and Bontemps and Meddahi (2012) for in-sample tests of distributional assumptions.

⁹Hong, Li and Zhao (2007) provide an out-of-sample counterpart of the Hong and Li (2005) in-sample tests, while Corradi and Swanson (2006b) generalize the in-sample test by Corradi and Swanson (2006a) to an out-of-sample framework.

our approach is complementary to theirs. Furthermore, since the null hypothesis is different, we cannot directly compare our tests to theirs. Knueppel (2015) tests the correct calibration of multi-step-ahead density forecasts using raw moments, that is, by testing whether selected moments of the PITs are the same as the corresponding moments of a uniform distribution. Berkowitz (2001) instead proposes a test based on the Inverse Normal transformation of the PITs. Both Berkowitz (2001) and Knueppel (2015) abstract from parameter estimation error.¹⁰ In Corollaries 5 and 6, we formally discuss the likelihood ratio and the raw-moments-based tests in our framework.

There are several differences between tests based on raw moments (such as Knueppel, 2015) and ours. On the one hand, it is important to note that raw-moments-based tests evaluate correct specification using a finite number of moments while our approach directly tests the correct specification of the whole distribution of the PITs. Therefore, Knueppel’s (2015) test has power to detect misspecification only if it includes the moments that capture misspecification, but would not have power if the misspecification affects moments that are not included. One of the drawbacks of Knueppel’s (2015) test, then, is that it requires the researcher to choose which moments to test and it is unclear how to select the number of moments to test. Our approach, instead, is equivalent to testing the correct calibration of the whole distribution, which corresponds to testing the correct specification of all the moments and does not suffer from this drawback. On the other hand, the fact that Knueppel’s (2015) test relies on a finite number of moments gives it two advantages: first, since the inclusion of additional, correctly-specified, moments comes with the cost of a power loss, it may have more power than our tests if the misspecification is fully captured by a few moments only and the researcher has chosen to test exactly those moments; second, in the case of serial correlation, Knueppel’s (2015) test relies on a small number of moments whose covariance can be consistently and precisely estimated using a HAC estimator, while in our case the covariance matrix is large dimensional and hence we recommend a bootstrap for a more precise and robust inference. The alternative approach recently proposed by González-Rivera and Sun (2014) uses graphical devices to implement a test of correct specification. Their proposed methods work when models are dynamically correctly specified; however, when parameter estimation error is asymptotically relevant, the asymptotic distribution is not nuisance parameter free and a bootstrap procedure is proposed. Our tests, instead, do not require a bootstrap procedure for one-step-ahead predictive densities, and its critical

¹⁰Knueppel (2015) conjectures that our framework can be applied to his approach. We formally show under which assumptions his conjecture is valid.

values are readily available.

To summarize, for one-step-ahead predictive densities, the critical values of our test statistics can be tabulated. For multiple-step-ahead forecasts, the PITs are potentially serially correlated; in the latter case, we recommend to obtain the critical values of our test statistics via a block version of the weighted bootstrap proposed by Inoue (2001). An alternative would be to obtain the critical values via Monte Carlo simulations using a HAC covariance estimate, but the performance is sensitive to the choice of the truncation lag and the correlation properties in the data, and thus it could perform poorly in practice. Even though we implement a bootstrap to obtain the critical values, our bootstrap is very different from Corradi and Swanson (2006a). To highlight the theoretical differences between the null hypothesis in our framework and Corradi and Swanson's (2006a), note that the variance in the limiting distribution of the test statistic in Corradi and Swanson (2006a) includes the contribution of parameter estimation error to the variance. Our proposed weighted block bootstrap is different than that in Corradi and Swanson (2006a) since it resamples the PITs and not the data, and does not require knowledge or replication of the model estimation technique. Moreover, as previously noted, critical values are readily available for one-step-ahead predictive densities, and a bootstrap procedure is only needed when evaluating multi-step predictive densities.

To illustrate the empirical relevance of our proposed tests, we evaluate density forecasts in the Survey of Professional Forecasters (SPF). It is very interesting to evaluate the SPF density forecasts using our tests because SPF panelists use a combination of estimated models and expert judgment to produce forecast, even though the models are not disclosed. Thus, in the SPF density forecast case, as well as in most central banks' fan charts, there is parameter estimation error, and it is impossible to correct for it: the only feasible approach is to maintain it under the null hypothesis. This example illustrates the empirical usefulness of our tests, since this approach is exactly what we follow in our paper. In fact, we find that predictive densities are, in general, misspecified. In addition, we propose ways to improve the calibration of the densities given the results of our tests.

The remainder of the paper is organized as follows. Section 2 introduces the notation and definitions, and Section 3 discusses issues related to the practical applicability of our test as well as our theoretical approach. In Section 4, we provide Monte Carlo evidence on the performance of our tests in small samples. Section 5 analyzes the empirical applications to SPF density forecasts, and Section 6 concludes.

2 Notation and Definitions

We first introduce the notation and discuss the assumptions about the data, the models and the estimation procedure. Consider a stochastic process $\{Z_t : \Omega \rightarrow \mathbb{R}^{k+1}\}_{t=1}^T$ defined on a complete probability space (Ω, \mathcal{F}, P) . The observed vector Z_t is partitioned as $Z_t = (y_t, X_t)'$, where $y_t : \Omega \rightarrow \mathbb{R}$ is the variable of interest and $X_t : \Omega \rightarrow \mathbb{R}^k$ is a vector of predictors. Let $1 \leq h < \infty$.¹¹ We are interested in the true, but unknown, h -step-ahead conditional predictive density for the scalar variable y_{t+h} based on $\mathcal{F}_t = \sigma(Z_1', \dots, Z_t')$, which is the true information set available at time t . We denote this density by $\phi_0(\cdot)$.¹²

We assume that the researcher has divided the available sample of size $T + h$ into an in-sample portion of size R and an out-of-sample portion of size P and obtained a sequence of h -step-ahead out-of-sample density forecasts of the variable of interest y_t using the information set \mathfrak{S}_t , such that $R + P - 1 + h = T + h$ and $\mathfrak{S}_t \subseteq \mathcal{F}_t$. Note that this implies that the researcher observes a subset of the true information set. We also let \mathfrak{S}_{t-R+1}^t denote the truncated information set between time $(t - R + 1)$ and time t used by the researcher.

Let the sequence of P out-of-sample estimates of conditional predictive densities evaluated at the ex-post realizations be denoted by $\{\phi_{t+h}(y_{t+h} | \mathfrak{S}_{t-R+1}^t)\}_{t=R}^T$. The dependence on the information set is a result of the assumptions we impose on the in-sample parameter estimates, $\hat{\theta}_{t,R}$. We assume that the parameters are re-estimated at each $t = R, \dots, T$ over a window of R data indexed from $t - R + 1$ to t (rolling scheme).¹³ In this paper we are concerned with direct multi-step forecasting, where the predictors are lagged h periods. In addition to being parametric (such as a Normal distribution), the distribution $\phi_{t+h}(\cdot)$ can also be non-parametric (as in one of the empirical applications in this paper).

Consider the probability integral transform (PIT), which is the cumulative density function (CDF) corresponding to $\phi_{t+h}(\cdot)$ evaluated at the realized value y_{t+h} :

¹¹Note that our framework allows nowcast densities, which technically corresponds to $h = 0$. We do not make this explicit in the notation to avoid misleading the reader to thinking that our tests are in-sample.

¹²The true conditional predictive density may depend on the forecast horizon. To simplify notation, we omit this dependence without loss of generality given that the forecast horizon is fixed. Furthermore, we use the symbols $\phi_0(\cdot)$ and $\phi_t(\cdot)$ to denote generic distributions and not necessarily a normal distribution.

¹³The choice of the estimation scheme (rolling versus recursive) depends on the features of the data: in the presence of breaks, one would favor a rolling scheme that allows a fast update of the parameter estimates, at the cost of a potential increase in estimation uncertainty relative to a recursive scheme when there are no breaks. As discussed in Giacomini and White (2006), our proposed approach is also valid for other classes of limited memory estimators.

$$z_{t+h} = \int_{-\infty}^{y_{t+h}} \phi_{t+h}(y|\mathfrak{S}_{t-R+1}^t) dy \equiv \Phi_{t+h}(y_{t+h}|\mathfrak{S}_{t-R+1}^t).$$

Let us also denote the empirical cumulative probability distribution function of the PIT by

$$\varphi_P(r) \equiv P^{-1} \sum_{t=R}^T \mathbf{1} \{ \Phi_{t+h}(y_{t+h}|\mathfrak{S}_{t-R+1}^t) \leq r \}. \quad (1)$$

Further, let

$$\xi_{t+h}(r) \equiv (1 \{ \Phi_{t+h}(y_{t+h}|\mathfrak{S}_{t-R+1}^t) \leq r \} - r),$$

where $1 \{.\}$ is the indicator function and $r \in [0, 1]$. Consider $\Psi(r) = \Pr \{ z_{t+h} \leq r \} - r$ and its (rescaled) out-of-sample counterpart:

$$\Psi_P(r) \equiv P^{-1/2} \sum_{t=R}^T \xi_{t+h}(r). \quad (2)$$

3 Asymptotic Tests of Correct Specification

This section presents our proposed test for the case of one-step-ahead forecasts. In this case, our tests have an asymptotic distribution that is free of nuisance parameters and the critical values can be tabulated. We further generalize the tests to the presence of serial correlation, which applies to the case of multi-step-ahead density forecasts. In this case the asymptotic distribution is not nuisance parameter free, and we discuss how to calculate the critical values. All the proofs are relegated to Appendix A.

In order to maintain parameter estimation error under the null hypothesis, we state our null hypothesis in terms of a truncated information set, which expresses the dependence of the predictive density on estimated parameter values (as in Amisano and Giacomini, 2007). We focus on testing $\phi_{t+h}(y|\mathfrak{S}_{t-R+1}^t) = \phi_0(y|\mathcal{F}_t)$, that is:

$$H_0 : \Phi_{t+h}(y|\mathfrak{S}_{t-R+1}^t) = \Phi_0(y|\mathcal{F}_t) \text{ for all } t = R, \dots, T, \quad (3)$$

where $\Phi_0(y|\mathcal{F}_t) \equiv \Pr(y_{t+h} \leq y|\mathcal{F}_t)$ denotes the distribution specified under the null hypothesis.¹⁴ The alternative hypothesis, H_A , is the negation of H_0 . Note that the null hypothesis

¹⁴Note that the null hypothesis depends on R . In other words, the null hypothesis jointly tests density functional form and estimation technique. It might be possible that correct specification is rejected for a model for some values of R and not rejected for the same model for some other choices of R . This is reasonable since we are evaluating the model's performance when estimated with a given sample size, so

evaluates the correct specification of the density forecast of a model estimated with a given window size, R , as well as the parameter estimation method chosen by the researcher.

We are interested in the test statistics:

$$\kappa_P = \sup_{r \in [0,1]} |\Psi_P(r)|, \quad (4)$$

$$C_P = \int_0^1 \Psi_P(r)^2 dr. \quad (5)$$

The κ_P test statistic is the same as the V_{1T} test statistic considered by Corradi and Swanson (2006a) when applied to predictive densities. Note, however, that we derive the asymptotic distribution of the test statistic under a different null hypothesis. Corradi and Swanson (2006a) focus on the null hypothesis: $H_0^{CS} : \Phi_{t+h}(y|\mathfrak{S}_t) = \Phi_0(y|\mathfrak{S}_t, \theta^\dagger)$ for some pseudo-true parameter value $\theta^\dagger \in \Theta$, where Θ is the parameter space. That is, the latter test the hypothesis of correct specification of the predictive density at the pseudo-true parameter value. Thus, the limiting distribution of their test reflects parameter estimation error and, therefore, is not nuisance parameter free. Note that we cannot compare our test with Corradi and Swanson (2006a) since they focus on a different null hypothesis where R tends to infinity, while the theory of our test relies on R being finite. In fact, given that the null hypotheses are different, power in our context corresponds to size in theirs; thus comparisons are not informative.

Under our null hypothesis (eq. 3) instead, the limiting distribution of the test statistic is nuisance parameter free. The reason is that we maintain parameter estimation error under the null hypothesis, which implies that the asymptotic distribution of the test does not require a delta-method approximation around the pseudo-true parameter value.

To clarify our null hypothesis, we provide a few examples.

Example 1: As a simple example, consider $y_{t+1} = c_{t+1} + x_t + \varepsilon_{t+1}$, $\varepsilon_{t+1} \sim iid N(0, 1)$, where *iid* means independent and identically distributed, $x_t \sim iid N(0, 1)$, $t = 1, \dots, T$ and ε_{t+1} , x_t are independent of each other. We assume for simplicity that the variance of the errors is known. The researcher instead considers a model $y_{t+1} = \beta x_t + e_{t+1}$, $e_{t+1} \sim iid N(0, 1)$. Moreover, the researcher is re-estimating the coefficient β with a window of size R at each point in time t . Let $\hat{\beta}_{t,R}$ denote the parameter estimated at time t using the most recent R observations. We set c_{t+1} such that our null hypothesis (eq. 3) holds. That is, the

the estimation error is important under the null hypothesis. Alternatively, one could construct a test that is robust to the choice of the estimation window size as suggested in Inoue and Rossi (2012) and references therein.

estimated PIT is:

$$\int_{-\infty}^{y_{t+1}} \phi_{t+1}(y|\mathfrak{S}_{t-R+1}^t) dy,$$

where $\phi_{t+1}(y|\mathfrak{S}_{t-R+1}^t)$ is $N(\widehat{\beta}_{t,R}x_t, 1)$, whereas the PIT of the true data generating process (DGP) is:

$$\int_{-\infty}^{y_{t+1}} \phi_{t+1}(y|\mathcal{F}_t) dy,$$

where $\phi_{t+1}(y|\mathcal{F}_t)$ is $N(c_{t+1} + x_t, 1)$. Under the assumption that the variance is known, a sufficient condition for the null hypothesis to hold is that the conditional means from true DGP and the estimated model are the same. More in detail, the null hypothesis is imposed by assuming:

$$c_{t+1} + x_t = \widehat{\beta}_{t,R}x_t,$$

that is,¹⁵

$$c_{t+1} = \left(\frac{R^{-1} \sum_{j=t-R+1}^t x_{j-1}y_j}{R^{-1} \sum_{j=t-R+1}^t x_{j-1}^2} - 1 \right) x_t.$$

Thus, the null hypothesis in eq. (3) is not the correct specification of the forecasting model evaluated at the true parameter values (relative to the data generating process); rather, the null hypothesis in eq. (3) is the correct specification of the forecasting model evaluated at the parameter values obtained conditional on the estimation procedure. We argue that the latter is an appropriate approach to evaluate the correct specification of density forecasts, since it jointly evaluates the proposed model and its estimation technique, including the estimation window size. The methodology only requires that the conditional mean is estimated based on a finite number of observations.¹⁶

Suppose, instead, the true data generating process is: $y_{t+1} = c + x_t + \varepsilon_{t+1}$ where $x_t \sim iid\chi_1^2$ and $\varepsilon_{t+1} \sim iidN(0, 1)$. Let the researcher estimate a misspecified model that includes only a constant, treating the forecast distribution as Normal. Note that the null hypothesis does not hold even if the error term is Normal, since the misspecification results in an actual error

¹⁵The data under the null hypothesis are mixing, and thus satisfy our Assumption 1, for the following reason: let $g_{t+1} \equiv (x_t, c_{t+1}, \varepsilon_{t+1})'$. Since $E(g_{t+1}) = 0$ and $E(g_{t+1}|g_t, g_{t-1}, \dots) = 0$ then g_{t+1} is a martingale difference sequence and has finite variance, thus it is white noise (Hayashi, 2000, p. 104).

¹⁶The results in this paper also carry over to the fixed-estimation scheme, where the conditioning information set is \mathfrak{S}_1^R , or to any other information set based on a bounded number of observations R , provided R is finite.

term that is a combination of x_t and ε_{t+1} . Thus, since the data is generated as a mixture of a chi-squared and Normal distribution, and we are testing whether it is a Normal, the null hypothesis does not hold.

Example 2: Consider $y_{t+1} = \sigma_{t+1}\varepsilon_{t+1}$, $\varepsilon_{t+1} \sim iid N(0, 1)$ and $\sigma_t^2 = \rho_1\sigma_{t-1}^2 + \rho_{0,t}$. This is a GARCH(1,0) process for y_{t+1} , where the mean is assumed to be known and equal to zero for simplicity. The researcher instead estimates the model: $y_{t+1} = \gamma e_{t+1}$, $e_{t+1} \sim iid N(0, 1)$, where the coefficient γ is estimated using observations in a window of size R : $\hat{\gamma}_t = R^{-1} \sum_{j=t-R+1}^t y_j^2$. That is, the estimated PIT is:

$$\int_{-\infty}^{y_{t+1}} \phi_{t+1}(y | \mathfrak{S}_{t-R+1}^t) dy,$$

where $\phi_{t+1}(y | \mathfrak{S}_{t-R+1}^t)$ is $N(0, \hat{\gamma}_t)$, whereas the PIT of the true data generating process (DGP) is:

$$\int_{-\infty}^{y_{t+1}} \phi_{t+1}(y | \mathcal{F}_t) dy,$$

where $\phi_{t+1}(y | \mathcal{F}_t)$ is $N(0, \sigma_{t+1}^2)$.

We set $\rho_{0,t}$ such that our null hypothesis (eq. 3) holds. The estimated PIT and the PIT that generated the data are the same if $\sigma_{t+1}^2 = \hat{\gamma}_t$. Thus, the null hypothesis is imposed by assuming:

$$\sigma_{t+1}^2 = R^{-1} \sum_{j=t-R+1}^t y_j^2 = R^{-1} \sum_{j=t-R+1}^t (\rho_1\sigma_{j-1}^2 + \rho_{0,j})^2 \varepsilon_j^2,$$

where $\sigma_{t+1}^2 = \rho_1\sigma_t^2 + \rho_{0,t+1}$ and

$$\rho_{0,t+1} = R^{-1} \sum_{j=t-R+1}^t y_j^2 - \rho_1\sigma_t^2 = R^{-1} \sum_{j=t-R+1}^t (\rho_1\sigma_{j-1}^2 + \rho_{0,j})^2 \varepsilon_j^2 - \rho_1\sigma_t^2.$$

Example 3: As an example of a dynamically misspecified model where the null hypothesis in eq. (3) holds, consider $y_{t+1} = c_{t+1} + \rho y_t + \varepsilon_{t+1}$, $\varepsilon_{t+1} \sim iid N(0, 1)$. We assume for simplicity that the variance of the errors is known. The researcher instead considers a model $y_{t+1} = \beta + e_{t+1}$, $e_{t+1} \sim iid N(0, 1)$. Moreover, the researcher is estimating the coefficient β using observations in a window of size R . That is, the estimated PIT is:

$$\int_{-\infty}^{y_{t+1}} \phi_{t+1}(y | \mathfrak{S}_{t-R+1}^t) dy,$$

where $\phi_{t+1}(y | \mathfrak{S}_{t-R+1}^t)$ is $N(\hat{\beta}_{t,R}, 1)$, where $\hat{\beta}_{t,R} = R^{-1} \sum_{j=t-R+1}^t y_j$, whereas the PIT of the true DGP is:

$$\int_{-\infty}^{y_{t+1}} \phi_0(y | \mathcal{F}_t) dy,$$

where $\phi_0(y|\mathcal{F}_t)$ is $N(c_{t+1} + \rho y_t, 1)$. Under the assumption that the variance is known, a sufficient condition for the null hypothesis to hold is that the conditional means from the true DGP and the estimated model are the same. More in detail, the null hypothesis is imposed by assuming:

$$c_{t+1} + \rho y_t = \widehat{\beta}_{t,R},$$

that is,

$$c_{t+1} = \left(R^{-1} \sum_{j=t-R+1}^t y_j - \rho y_t \right).$$

It is important to note that R is finite; thus $R^{-1} \sum_{j=t-R+1}^t y_j$ is a mixing process since it is a measurable function of a finite number of lags of mixing random variables. In summary, in this case, even if the forecasting model is misspecified relative to the data generating process, the null hypothesis in eq. (3), which aims to evaluate correct specification of the model and the forecasting technique jointly, holds.

3.1 One-step-ahead Density Forecasts

This sub-section presents results for the case of one-step-ahead forecasts; the next sub-section generalizes the tests to the presence of serial correlation. Let $h = 1$. First, we derive the asymptotic distribution of $\Psi_P(r)$ for one-step-ahead density forecasts under Assumption 1.¹⁷

Assumption 1.

(i) $\{Z_t = (y_t, X_t)'\}_{t=R}^T$ is strong mixing with mixing coefficients $\alpha(m)$ of size $-\delta$, where $\delta > 3(4 + \gamma)/\gamma$;

(ii) $\Phi_0(y_{t+h}|\mathcal{F}_t)$ is differentiable and has a well-defined inverse;

(iii) $F_d(\cdot, \cdot)$ and $F(\cdot)$ are respectively the joint and the marginal distribution functions of the random variable $\Phi_0(y_{t+h}|\mathcal{F}_t)$, i.e. $\Pr(\Phi_0(y_{t+h}|\mathcal{F}_t) \leq r_1, \Phi_0(y_{t+h+d}|\mathcal{F}_{t+d}) \leq r_2) = F_d(r_1, r_2)$, $\Pr(\Phi_0(y_{t+h}|\mathcal{F}_t) \leq r) = F(r)$, and $F(r)$ is continuous;

(iv) $R < \infty$ as $P, T \rightarrow \infty$.

Assumption 1(i) allows for short memory and heterogeneous data. The assumption however limits the dependence in the data so that, as in Giacomini and White (2006), one can

¹⁷Note that if $P/R \rightarrow 0$, our test would be the same as the existing tests as parameter estimation uncertainty becomes irrelevant in those cases (see Corradi and Swanson, 2006b). This result would hold even for recursive estimation schemes as long as $P/R \rightarrow 0$. However, we test a different null hypothesis than the existing tests and we do not allow either $R \rightarrow \infty$ or $P/R \rightarrow 0$ in our framework.

use results on functions of mixing variables which allow for mild non-stationarity induced by changes in distributions over time, yet rule out I(1) processes. Assumption 1(ii) and 1(iii) are similar to Assumption B in Inoue (2001); assumption (iii) is trivially satisfied under our null hypothesis since, as we will show, the PITs are uniformly distributed under the null hypothesis. These assumptions require the PITs, as well as the marginal and joint distributions of the PITs, to be well-defined.¹⁸ Assumption 1(iv) requires the estimation window size to be finite as the total sample size grows. Our assumptions allow for quite general parametric models (including linear and nonlinear models) and estimation methods (including GMM and MLE), as long as the estimation window size is finite and data are mixing. Note that the parameters of the model do not need to be consistently estimated as long as assumptions 1(i) and 1(iv) hold. Furthermore, note that the assumption potentially allows forecasts to be conditioned on a finite set of future values of some variables of interest (i.e. “conditional forecasts”).

Correct specification is characterized by Assumption 2:

Assumption 2.

$y_{t+h} | \mathfrak{S}_{t-R+1}^t \equiv y_{t+h} | \mathcal{F}_t$ for all $t = R, \dots, T$, where \equiv denotes equality in distribution.

We show the following result:

Theorem 1 (Asymptotic Distribution of $\Psi_P(r)$) Under Assumptions 1, 2, and H_0 in eq. (3): (i) z_{t+h} is iid $U(0, 1)$, $t = R, \dots, T$; (ii) $\Psi_P(r)$ weakly converges (\Rightarrow) to the Gaussian process $\Psi(\cdot)$, with mean zero and auto-covariance function $E[\Psi(r_1)\Psi(r_2)] = [\inf(r_1, r_2) - r_1r_2]$.

The result in Theorem 1 allows us to derive the asymptotic distribution of the test statistics of interest, presented in Theorem 2. The latter shows that the asymptotic distributions of our proposed test statistics have the appealing feature of being nuisance parameter free. Note that the asymptotic distribution is a Brownian Bridge (see Durbin, 1973).

Theorem 2 (Correct Specification Tests) Under Assumptions 1, 2 and H_0 in eq. (3):

$$\kappa_P \equiv \sup_{r \in [0,1]} |\Psi_P(r)| \Rightarrow \sup_{r \in [0,1]} |\Psi(r)|, \quad (6)$$

and

$$C_P \equiv \int_0^1 \Psi_P(r)^2 dr \Rightarrow \int_0^1 \Psi(r)^2 dr. \quad (7)$$

¹⁸The assumption is on the unobserved true distribution, though under the null it also ensures that the proposed distribution has a well defined limiting distribution.

The tests reject H_0 at the $\alpha \cdot 100\%$ significance level if $\kappa_P > \kappa_\alpha$ and $C_P > C_\alpha$. Critical values for $\alpha = 0.10, 0.05$ and 0.01 are provided in Table 1, Panel A.

INSERT TABLE 1 HERE

Note that one could be interested in testing correct specification in specific parts of the distribution.¹⁹ For example, one might be interested in the tails of the distribution, which correspond to outliers, such as the left tail, where $r \in [0, 0.25)$, or the right tail, where $r \in [0.75, 1)$, or both: $r \in \{[0, 0.25] \cup [0.75, 1]\}$. Alternatively, one might be interested in the central part of the distribution, for example $r \in [0.25, 0.75]$. We provide critical values for these interesting cases in Table 1, Panel B.

Note also that our κ_P test has a graphical interpretation. In fact, from eqs. (1) and (2),

$$P^{-1/2}\Psi_P(r) \equiv P^{-1} \sum_{t=R}^T (1 \{ \Phi_{t+h}(y_{t+h} | \mathcal{S}_{t-R+1}^t) \leq r \} - r) = \varphi_P(r) - r.$$

Thus,

$$\alpha \equiv \Pr \left\{ \sup_{r \in [0,1]} |\Psi_P(r)| > \kappa_\alpha \right\} \approx \Pr \left\{ \sup_{r \in [0,1]} |\varphi_P(r) - r| > \kappa_\alpha / \sqrt{P} \right\}.$$

This suggests the following implementation: plot the empirical cumulative distribution function of the PIT, eq. (1), together with the cumulative distribution function of the Uniform (0,1) distribution, r (the 45-degree line), and the critical value lines: $r \pm \kappa_\alpha / \sqrt{P}$. Then, the κ_P test rejects if the cumulative distribution function of the PIT is outside the critical value bands.

We consider two ways of simulating the critical values. One approach, which is what we report in Table 1, relies on simulating the limiting distribution of $\Psi_P(r)$, considered in Theorem 1, directly. More specifically,

- (i) Discretize the grid for r . In particular, we consider the grid: $\underline{r} = [0 : 0.001 : 1]$;
- (ii) Calculate the theoretical variance $E[\Psi(r_1)\Psi(r_2)] = [\inf(r_1, r_2) - r_1 r_2]$, for $r_1, r_2 \in [0 : 0.001 : 1]$;
- (iii) Draw independent multivariate Normal random variables based on the Cholesky decomposition of the estimated covariance matrix $E[\Psi(r_1)\Psi(r_2)]$ calculated in (ii);²⁰
- (iv) Construct the test statistics proposed in Theorem 2;

¹⁹See Franses and van Dijk (2003), Amisano and Giacomini (2007) and Diks, Panchenkob and van Dijk (2011) for a similar idea in the context of point and density forecast comparisons.

²⁰A finer grid, in theory, could result in a more accurate approximation, but numerically the Cholesky factorization becomes less accurate. Thus, there is a limited payoff from refining the grid.

(v) Repeat the steps (iii) and (iv) for a large number of Monte Carlo replications; report the 90%, 95% and 99% percentiles of the simulated limiting distribution as critical values.

The second approach aims at obtaining exact critical values in finite samples. We do so by the following procedure:

(i) Draw P independent random variables from the Uniform (0,1) distribution;

(ii) For a given sample of size P , construct the $\Psi_P(r)$ as in eq. (1);

(iii) Construct the test statistics proposed in Theorem 2;

(iv) Repeat steps (i) to (iii) for a large number of Monte Carlo replications; the 90%, 95% and 99% percentile values of the simulated limiting distribution are the critical values.

The second approach has the advantage that we can tailor the critical values to the sample sizes used in an empirical application. However, it has the disadvantage that the critical values need to be simulated for each empirical application. As we show later in a Monte Carlo size experiment (presented in Table 2), the finite-sample critical values improve the size, more so for smaller values of P .

In principle, the limiting distribution of the Kolmogorov-Smirnov test reported in Theorem 1 is not different from that of the usual textbook version of the Kolmogorov-Smirnov test that are derived by analytical calculations (see Durbin, 1973). According to Smirnov (1948) these values are 1.22, 1.36 and 1.63 for $\alpha = 0.10, 0.05$ and 0.01 , respectively. Unreported Monte Carlo size experiments suggest that our simulated critical values result in tests that have better size than those based on the theoretical formula: the theoretical formula appears to be more conservative than our simulation-based procedure.

It is interesting to compare our approach to Diebold et al. (1998). While our null hypothesis is different from theirs, the procedure that we propose is similar to theirs in that both their implementation and ours abstract from parameter estimation error (for different reasons). Thus, our approach can be viewed as a formalization of their approach, albeit with a different null hypothesis. An additional advantage of our approach is that the critical value bands that we propose are joint, not pointwise.

The previous discussion suggests that we could also apply our approach to likelihood-ratio (LR) tests based on the Inverse Normal transformation of the PITs or raw-moment-based tests on the PITs. We will discuss such approaches in the next section.

Finally, note that our approach provides not only a rationale to the common practice of evaluating the correct specification of density forecasts using PITs without adjusting for parameter estimation error (Diebold et al., 1998), but also a methodology for implementing tests robust to the presence of serial correlation. This more general case is considered next.

3.2 Multi-step-ahead Forecasts

When considering h-step-ahead forecasts, $h > 1$ and finite, an additional problem arises as the PITs become serially correlated. Thus, we need to extend our results and allow the forecasts to be serially correlated under the null hypothesis; that is, when Assumption 2 does not hold.

When evaluating h-step-ahead conditional predictive densities, the next Theorem shows that $\Psi_P(r)$ weakly converges to the Gaussian process $\Psi(\cdot, \cdot)$, with mean zero and an auto-covariance function that depends on the serial correlation in the PITs.

Theorem 3 (Correct Specification Tests under Serial Correlation) *Under Assumption 1 and H_0 in eq. (3): (i) z_{t+h} is $U(0, 1)$, $t = R, \dots, T$; (ii) $\Psi_P(r)$ weakly converges to the Gaussian process $\Psi(\cdot, \cdot)$, with mean zero and auto-covariance function $E[\Psi(r_1)\Psi(r_2)] = \sigma(r_1, r_2)$, where $\sigma(r_1, r_2) = \sum_{d=-\infty}^{\infty} [F_d(r_1, r_2) - F(r_1)F(r_2)]$. (iii) Furthermore,*

$$\kappa_P \Rightarrow \sup_{r \in [0,1]} |\Psi(r)|,$$

$$C_P \Rightarrow \int_0^1 \Psi(r)^2 dr.$$

In this case, the limiting distribution resembles the one that Corradi and Swanson (2006c) obtain under dynamic misspecification since the limiting distribution is not free from parameter estimation error. However, under the null hypothesis, we are not concerned about dynamic misspecification since the null hypothesis may hold even though a model can be dynamically misspecified (see Example 3 in the previous section). Furthermore, to highlight the theoretical differences between the null hypothesis in our framework and Corradi and Swanson's (2006a), note that the variance in the limiting distribution of the test statistic in Corradi and Swanson (2006a) is more complicated since it includes the contribution of parameter estimation error to the variance.

Theorem 3 shows that, in the presence of serial correlation, the critical values depend on nuisance parameters that appear in the covariance matrix of the PITs. Let $r \in [0, 1]$ be discretized over a p -dimensional grid, $\underline{r} \equiv [r_1, r_2, \dots, r_p]'$, where p is large. Although Inoue (2001) conjectures that it might be possible to consistently estimate the $(p \times p)$ dimensional covariance of $[\xi_{t+h}(r_1), \xi_{t+h}(r_2), \dots, \xi_{t+h}(r_p)]'$ using a standard Newey and West's (1987) HAC estimator, unreported simulations show that, in practice, the estimator is sensitive to the choice of the bandwidth and the serial correlation properties of the data (see also Corradi, Jin and Swanson, 2016). We therefore recommend using critical values of κ_P and C_P from

a block version of the weighted bootstrap proposed by Inoue (2001). The assumptions and the implementation of the bootstrap are described in detail in what follows.

Assumption 1(i’). Let $Q \geq 16$ be an even integer and $\{Z_t = (y_t, X_t)'\}_{t=R}^T$ be strong mixing with mixing coefficients $\alpha(m)$ of size $-\delta$, where $\delta > (1 + Q/2) \frac{Q+\gamma}{\gamma}$.

Assumption 3. Let ℓ be the block length and η_t be a continuous random variable that is used for random weighting in the block weighted bootstrap. $\{\eta_t\}_{t=R}^{T-\ell+1}$ are independent random variables, independent of $\{z_{t+h}\}$, with zero mean, variance $1/\ell$ and $E(\eta_t^4) = O(1/\ell^2)$, where $\ell \rightarrow \infty$ as $T \rightarrow \infty$ and $\ell = o(P^{1/2})$.

Assumption 3 is the same as Assumption D in Inoue (2001). We follow Inoue (2001) in using $\eta_t \sim iidN(1, 1/\ell)$ in practice when constructing the bootstrap statistics. Let ω be a particular bootstrap sample. Let $z_{t+h}(\omega)$ be a realization of the PIT in a particular bootstrap sample. Note that the bootstrap we describe below requires resampling the PITs, z_{t+h} , not the data Z_t : this is the fundamental difference between our bootstrap approach and Corradi and Swanson’s (2006a). Define the bootstrap test statistics as $\Psi_P^*(r; \omega) = P^{-1/2} \sum_{j=R}^{T-\ell+1} \eta_j \sum_{i=j}^{j+\ell-1} (1 \{z_{t+h}(\omega) \leq r\} - r)$ and $\kappa_P^*(\omega) = \sup_{r \in [0,1]} |\Psi_P^*(r; \omega)|$, $C_P^*(\omega) = \int_0^1 \Psi_P^*(r; \omega)^2 dr$.

Theorem 4 (Bootstrap Validity) Under Assumptions 1(i’,ii-iv), 3 and H_0 in eq. (3):
(i) $\Psi_P^*(\cdot, \omega) \Rightarrow \Psi(\cdot)$; (ii) $\kappa_P^*(\omega) \Rightarrow \sup_{r \in [0,1]} |\Psi(r)|$; and (iii) $C_P^*(\omega) \Rightarrow \int_0^1 \Psi(r)^2 dr$, $\omega - a.s.$

The bootstrap can be implemented in practice using the following step-by-step procedure:

- (i) Construct the test statistics as in Theorem 2;
- (ii) Let S be the maximum number of bootstrap replications. For $s = 1, 2, \dots, S$, generate $\{\kappa_{P;s}^*\}_{s=1}^S$ and $\{C_{P;s}^*\}_{s=1}^S$, where $\kappa_{P;s}^*$ and $C_{P;s}^*$ are based on $\left\{ \eta_t^{(s)} \right\}_{t=R}^{T-\ell+1}$;
- (iii) Estimate the level- α critical values $\widehat{c}_{\kappa,\alpha}^S$ and $\widehat{c}_{C,\alpha}^S$ from $\{\kappa_{P;s}^*\}_{s=1}^S$ and $\{C_{P;s}^*\}_{s=1}^S$, respectively.

There are several alternative solutions proposed in the literature that one could use within our approach as well. One approach is to discard data by reducing the effective sampling rate to ensure an uncorrelated sample (Persson, 1974 and Weiss, 1973). If the PITs are $(h - 1)$ -dependent, this can be implemented in practice by creating sub-samples of predictive distributions that are h periods apart. However, this procedure may not be possible in small samples, since the sub-samples may significantly reduce the size of the sample. In those cases, one may implement the procedure in several uncorrelated sub-samples of forecasts that are at least h periods apart and then use Bonferroni methods to obtain a joint test

without discarding observations (see Diebold et al., 1998). However, it is well-known that Bonferroni methods are conservative; thus the latter procedure, while easy to implement, may suffer from low power.

Note that our approach can be used to implement likelihood-ratio (LR) tests based on the Inverse Normal transformation of the PITs or Wald-type raw-moment-based tests on the PITs. As noted in the literature, these approaches have typically abstracted from parameter estimation uncertainty. When focusing on the traditional null hypothesis, H_0^{CS} , ignoring parameter estimation error leads to size distortions. Note that the size distortion is not only a small sample phenomenon, but persists asymptotically. The next result shows that, since parameter estimation error is maintained under our null hypothesis H_0 , eq. (3), there is no need to correct the asymptotic distribution and the implied critical values of Berkowitz's (2001) likelihood ratio tests and Knueppel's (2015) raw-moment-based to account for parameter estimation error.

Assumption 1(i''). $\{Z_t = (y_t, X_t')'\}_{t=R}^T$ is strong mixing with mixing coefficients $\alpha(m)$ of size $-\lambda/(\lambda - 2)$, where $\lambda > 2$.

Assumption 4. Let $s = s_1, s_2, \dots, s_N$, where N is finite. Let $\mathcal{H}(\cdot)$ be a real-valued function and $\underline{\mathcal{H}}_{t+h} \equiv [\mathcal{H}(z_{t+h})^{s_1}, \mathcal{H}(z_{t+h})^{s_2}, \dots, \mathcal{H}(z_{t+h})^{s_N}]$.

(i) $E|\mathcal{H}(z_{t+h})^{s}|^{2(\lambda+\varkappa)} < \infty$ for $\varkappa > 0$ for all t ;

(ii) $\Omega_P \equiv P^{-1} \sum_{t=R}^T E(\underline{\mathcal{H}}_{t+h} \underline{\mathcal{H}}'_{t+h}) + P^{-1} \sum_{j=1}^{h-1} \sum_{t=R+j}^T [E(\underline{\mathcal{H}}_{t+h} \underline{\mathcal{H}}'_{t+h-j}) + E(\underline{\mathcal{H}}_{t+h-j} \underline{\mathcal{H}}'_{t+h})]$ is uniformly positive definite.

Corollary 5 (Inverse Normal Tests) Let $\oplus^{-1}(\cdot)$ denote the inverse of the standard Normal distribution function. Under Assumptions 1(i''), 4 (holding for $s = s_1, s_2$) and H_0 in eq. (3): $\mathcal{H}(z_{t+h}) \equiv \oplus^{-1}(z_{t+h})$ is $N(0, 1)$.²¹

Thus, one could test for the correct specification of the density forecast by testing the correct specification of specific moments of $\oplus^{-1}(z_{t+h})$. For example, the researcher could estimate an AR(1) model for $\oplus^{-1}(z_{t+h})$ and test that the mean and the slope are both zero, and that the variance is one. This approach has the advantage of being informative regarding the possible causes underlying the misspecification of the density forecast, as it can focus on different moments, and it may perform better in small samples. The disadvantage of the approach is that, unlike the κ_P and C_P tests, it focuses on specific moments of the distribution rather than the whole (non-parametric) cumulative distribution function.

²¹Under Assumption 2, Corollary 5 implies $\mathcal{H}(z_{t+h}) \equiv \Phi^{-1}(z_{t+h})$ is iid $N(0, 1)$.

An alternative test has been proposed by Knueppel (2015). Knueppel (2015) tests the correct calibration of multi-step-ahead density forecasts using raw moments. In Corollary 6 below, we formalize Knueppel’s (2015) test in our framework.²²

Corollary 6 (Raw-Moments-Based Tests) *Let $m_s \equiv E[\mathcal{H}(z_{t+h})^s]$ be the moments of real-valued function transformation of the PITs. Also let $\hat{m}_s = P^{-1} \sum_{t=R}^T \mathcal{H}(z_{t+h})^s$ denote the estimated moments, $\hat{D}_P = [\hat{m}_{s_1} - m_{s_1}, \hat{m}_{s_2} - m_{s_2}, \dots, \hat{m}_{s_N} - m_{s_N}]'$ and $\hat{\Omega}_P$ is a consistent estimate of Ω_P . Under Assumptions 1(i), ii, iv) and H_0 in eq. (3): $\alpha_P \equiv P \hat{D}_P \hat{\Omega}_P^{-1} \hat{D}_P \rightarrow \chi_N^2$.*

There are several differences between the κ_P, C_P tests and the α_P test. On the one hand, it is important to note that raw-moment-based tests only test a finite number of moments (e.g. the mean and the variance of the PITs should equal the mean and the variance of a Uniform distribution) while our approach directly tests the whole distribution of the PITs. Therefore, the α_P test has power to detect misspecification only if it includes the moments that capture misspecification, but would not have power if the misspecification affects moments that are not included in \mathcal{H}_{t+h} . One of the drawbacks of the α_P test, then, is that it requires the researcher to choose which moments to test and it is unclear how to select the number of moments to test. Our approach, instead, is equivalent to testing the whole distribution, which corresponds to testing all the moments. Thus, the κ_P, C_P tests instead do not suffer from this drawback. On the other hand, the fact that the α_P test relies on a finite number of moments gives it two advantages: first, it may have more power than our test if the misspecification is fully captured by a few moments only and the researcher has chosen to test exactly those moments; second, that it is possible to consistently estimate the covariance matrix using a HAC estimator in the case of serial correlation, while in the κ_P, C_P tests the covariance matrix is large dimensional: the latter might be a concern for HAC estimation (see Corradi, Jin and Swanson, 2016) and hence we recommend a bootstrap procedure in the case of serial correlation. Furthermore, the α_P test has a limiting distribution that is nuisance parameter free, while the κ_P, C_P tests have a limiting distribution that is nuisance parameter free only for one-step-ahead forecasts.

4 Monte Carlo Evidence

In this section we analyze the size and power properties of our proposed tests in small samples for both one- and multi-step-ahead forecasting models. Note that comparisons with

²²We thank a referee for suggesting to include this discussion.

alternative methods (such as Corradi and Swanson, 2006c, or González-Rivera and Yoldas, 2012) are not meaningful since we focus on a null hypothesis that is different from theirs.

4.1 Size Analysis

To investigate the size properties of our tests we consider several Data Generating Processes (DGPs). The forecasts are based on model parameters estimated in rolling windows for $t = R, \dots, T + h$. We consider several values for in-sample estimation window of $R = [25, 50, 100, 200]$ and out-of-sample evaluation period $P = [25, 50, 100, 200, 500, 1000]$ to evaluate the performance of the proposed procedure. While our Assumptions require R finite, we consider both small and large values of R to investigate the robustness of our methodology when R is large.²³ The DGPs are the following:

DGP S1 (Baseline Model): We estimate a model $y_t = \beta x_{t-1} + e_t$, $e_t \sim iidN(0, 1)$. The data is generated by $y_t = \mu_t + x_{t-1} + \varepsilon_t$, $\varepsilon_t \sim iid N(0, 1)$ and $x_t \sim iid N(0, 1)$, where

$$\mu_t = \left(\frac{R^{-1} \sum_{j=t-R}^{t-1} x_{j-1} y_j}{R^{-1} \sum_{j=t-R}^{t-1} x_{j-1}^2} - 1 \right) x_{t-1}.$$

DGP S2 (Extended Model): We parameterize the model according to the realistic situation where the researcher is interested in forecasting one-quarter-ahead U.S. real GDP growth with the lagged term spread from 1959:I-2010:III. We estimate a model $y_t = \beta x_{t-1} + e_t$, $e_t \sim iidN(0, 1)$, while the data has been generated with the DGP: $y_t = \mu_t + \gamma x_{t-1} + \varepsilon_t$, $\varepsilon_t \sim iidN(0, 1)$, $x_t = 0.2 + 0.8x_{t-1} + \nu_t$, $\nu_t \sim iid N(0, 1.08^2)$ and is independent from ε_t , $\gamma = 0.48$ and

$$\mu_t = \left(\frac{R^{-1} \sum_{j=t-R+1}^{t-1} x_{j-1} y_j}{R^{-1} \sum_{j=t-R}^{t-1} x_{j-1}^2} - \gamma \right) x_{t-1}.$$

DGP S3 (GARCH): Consider the data being generated by a GARCH(1,0), where $y_t = \sigma_t \varepsilon_t$, $\varepsilon_t \sim iidN(0, 1)$ and the $\sigma_t^2 = \rho_{0,t} + \rho_1 \sigma_{t-1}^2$. On the other hand, the forecasting model is

²³Note that Corradi and Swanson (2006b) focus on a different null hypothesis based on $R \rightarrow \infty$, and the theory of our test instead relies on R being finite. In fact, given that the null hypotheses are different, power in our context corresponds to size in their context; thus comparisons are meaningless.

$y_t = \gamma_t e_t$, $e_t \sim iidN(0, 1)$, and

$$\rho_{0,t} = \left(R^{-1} \sum_{j=t-R}^{t-1} (\rho_{0,j} + \rho_1 \sigma_{j-1}^2) \varepsilon_j^2 - \rho_1 \sigma_{t-1}^2 \right).$$

DGPs S1-S3 are based on one-step-ahead forecast densities. DGP S4 considers the case of h -step-ahead forecast densities ($h = 2$) where the PITs are serially correlated by construction.

DGP S4 (Serial Correlation): The DGP is $y_t = \mu_t + x_{t-1} + \varepsilon_t + \rho \varepsilon_{t-1}$, $\varepsilon_t \sim iidN(0, 1)$, $x_t \sim iid N(0, 1)$, $\rho = 0.2$ and μ_t is as defined in DGP S1. The estimated model is: $y_t = \beta x_{t-1} + e_t$, $e_t \sim iid N(0, 1 + \rho^2)$.

DGP S5 (IMA Model): The DGP is $\Delta y_t = \mu_t + \varepsilon_t - \rho \varepsilon_{t-1}$, $\varepsilon_t \sim iidN(0, 1.261)$, $\rho = 0.275$ and μ_t is defined as

$$\mu_t = R^{-1} \sum_{j=t-R}^{t-1} \Delta y_j.$$

The parameters for the Monte Carlo design are from Stock and Watson (2007, Table 3); we pick their parameterization for the 1960:I-1983:IV sample period, i.e. the period of Great Inflation, a period when there is more variability in inflation. The estimated model is: $\Delta y_t = \beta + e_t$, $e_t \sim iid N(0, 1 + \rho^2)$.

We also consider a modified version of the IMA model to understand the behavior of our tests for multi-step-ahead forecast horizons. In total, we consider three data generating processes, $\Delta y_t = \mu_t + \varepsilon_t - \sum_{j=1}^p \rho^j \varepsilon_{t-j}$, where $p = 1, 3, 11$, corresponding to two-, four- and twelve-step-ahead forecast horizons.

The results for DGP S1 are shown in Table 2. The table shows that our proposed tests have good size properties. Furthermore, the finite sample critical values improve the test performance for small values of P (see Panel B). Table 3 shows that our tests perform well in finite samples in DGPs S2-S5, except for the smallest value of P ($P = 25$) in the case of multi-step-ahead forecasts. The comparison of the last two panels (Panels E and F) in Table 3 show that the size of the tests is robust to alternative block lengths used in the bootstrap procedure – again, except for the smallest sample size, $P = 25$.

INSERT TABLES 2 AND 3 HERE

4.2 Power Analysis

To investigate the power properties of our tests, we consider the case of misspecification in the following DGPs.

DGP P1: The data are generated from a linear combination of Normal and χ_1^2 distributions: $y_t = \mu_t + x_{t-1} + (1 - c)\hat{\sigma}_t\eta_{1,t} + c(\eta_{2,t}^2 - 1)/\sqrt{2}$, where x_t , $\eta_{1,t}$ and $\eta_{2,t}$ are *iidN*(0, 1) random variables that are independent of each other and μ_t is as defined in DGP S1. The researcher tests whether the data result from a Normal distribution, i.e. considers the model $y_t = \beta x_{t-1} + e_t$, $e_t \sim iidN(0, \sigma_e)$. When c is zero, the null hypothesis is satisfied. When c is positive, the density becomes a convolution of a standard Normal and a χ_1^2 distribution (with mean zero and variance one), where the weight on the latter becomes larger as c increases.²⁴

DGP P2: We estimate a model $y_t = \beta x_{t-1} + e_t$, $e_t \sim iidN(0, 1)$. The data is generated by $y_t = \mu_t + x_{t-1} + \varepsilon_t$, $\varepsilon_t \sim iid t_\nu$, where $x_t \sim iidN(0, 1)$, μ_t is defined as in DGP S1, while ν is the number of degrees of freedom. When ν is large, the null hypothesis is satisfied; as ν decreases, the misspecification increases.

The results shown in Table 4 suggest that our proposed specification tests (κ_P, C_P) have good power properties in detecting misspecification in the predictive density.²⁵

INSERT TABLE 4 HERE

5 Empirical Analysis

This section provides an empirical assessment of the correct specification of the Survey of Professional Forecasters (SPF) density forecasts of inflation and output growth. The reasons why we focus on this example is as follows. SPF panelists use a combination of estimated models and expert judgment to forecast, even though the models are not known and, even if they were, they apply expert judgment to trim the forecasts and/or combine models' forecasts. In fact, in a recent SPF survey overview, Stark (2013, p. 2) found that: "Overwhelmingly, the panelists reported using mathematical models to form their projections. However, we also found that the panelists apply subjective adjustments to their pure-model forecasts. The relative role of mathematical models changes with the forecast

²⁴Note that $(\eta_{2,t}^2 - 1)/\sqrt{2}$ is a chi-squared distribution with zero mean and variance one, that is, it has the same mean and variance as the normal distribution we have under the null hypothesis, although the shape is different.

²⁵Unreported results show that the test still has power when we consider smaller sample sizes, e.g. $T = 100$.

horizon.” Interestingly, the survey also found that SPF panelists “change their forecasting approach with the length of the forecast horizon. At the shortest horizons (two years out and less), mathematical models are widely used by the panelists. Between 18 to 20 forecasters reported using models at these short horizons (...). Panelists also reported using models for long-horizon projections as well (three or more years out), although proportionately fewer rely on models at the long horizons than at the short horizons. (...) They use a combination of models in forming their expectations, rather than just one model.” Thus, in the SPF density forecast case, it is impossible to correct for parameter estimation error: the only approach is to maintain it under the null hypothesis. This is exactly the approach we follow in our paper. This highlights the empirical usefulness of the methodologies described in our paper. In fact, one of the advantages of our testing approach is that the only information needed for the implementation is a predictive density: knowledge of the model that generated the forecasts is not necessary.

Diebold et al. (1999) evaluate the correct specification of the density forecasts of inflation in the SPF.²⁶ In this section, we conduct a formal test of correct specification for the SPF density forecasts using our proposed procedure and compare our results to theirs. In addition to inflation, we also investigate the conditional density forecasts of output growth.

We use real GNP/GDP and the GNP/GDP deflator as measures of output and prices. The mean probability distribution forecasts are obtained from the Federal Reserve Bank of Philadelphia (Croushore and Stark, 2001). In the SPF data set, forecasters are asked to assign a probability value (over pre-defined intervals) of year-over-year inflation and output growth for the current (nowcast) and following (one-year-ahead) calendar years. The forecasters update the assigned probabilities for the nowcasts and the one-year-ahead forecasts on a quarterly basis. The probability distribution provided by the SPF is discrete, and we base our results on a continuous approximation by fitting a Normal distribution. The realized values of inflation and output growth are based on the real-time data set for macroeconomists, also available from the Federal Reserve Bank of Philadelphia.

The analysis of the SPF probability distribution is complicated since the SPF questionnaire has changed over time in various dimensions: there have been changes in the definition of the variables, the intervals over which probabilities have been assigned, as well as the time

²⁶The SPF provides two types of density forecasts: one is the distribution of point forecasts across forecasters (which measures the dispersion of point forecasts across forecasters), and the other is the mean of the probability density forecasts (which measures the average of the density forecasts across forecasters). We focus on the latter.

horizon for which forecasts have been made. To mitigate the impact of these problematic issues, we truncate the data set and consider only the period 1981:III-2011:IV. To evaluate the density forecasts we use the year-over-year growth rates of output and prices calculated from the first quarterly vintage of the real GNP/GDP and the GNP/GDP deflator in levels. For instance, in order to obtain the growth rate of real output for 1981, we take the 1982:I vintage of data and calculate the growth rate of the annual average GNP/GDP from 1980 to 1981. We consider the annual-average over annual-average percent change (as opposed to fourth-quarter over fourth-quarter percent change) in output and prices to be consistent with the definition of the variables that SPF forecasters provide probabilistic predictions for.

The empirical results are shown in Table 5. Asterisks (“*”) indicate rejection at the 5% significance level based on the critical values in Theorem 2 (reported in Table 1, Panel A). The test rejects correct specification for both output growth and inflation, except for output growth at the one-year-ahead forecast horizon.

INSERT TABLE 5 HERE

Our results are important in light of the finding that survey forecasts are reportedly providing the best forecasts of inflation. For example, Ang et al. (2007) find that survey forecasts outperform other forecasting methods (including the Phillips curve, the term structure and ARIMA models) and that, when combining forecasts, the data put the highest weight on survey information. Our results imply that, in contrast, survey forecasts do not characterize the predictive distribution of inflation correctly.

Figure 2 plots the empirical CDF of the PITs (solid line). Under the null hypothesis in Theorem 2, the PITs should be uniformly distributed; thus the CDF of the PITs should be the 45 degree line. The figure also reports the critical values based on the κ_P test. If the empirical CDF of the PITs is outside the critical value bands, we conclude that the density forecast is misspecified. Clearly, the correct specification is rejected in all cases except the one-year-ahead density forecast of GDP growth.

The figure also provides a visual analysis of the misspecification in the PITs. For instance, in the case of the current year output growth forecasts in Panel A of Figure 2, it appears that there are not as many realization in the left tail of the distribution relative to what the forecasters expected (the slope in the left tail is flat relative to the 45 degree line). In the case of the current year inflation forecast in Panel B, the forecasters overestimate both tails of the distribution and, instead, do not put as much probability on potential outcomes in the middle of the distribution. One-year-ahead inflation forecasts appear to have a different

dynamics: there is no evidence of misspecification in the left tail, but there are many more realizations relative to expected frequencies in the left half of the distribution. This also comes at the expense of misspecifying the right tail of the distribution: the forecasters are more optimistic about extremely positive output growth scenarios relative to the realized outcomes.

INSERT FIGURES 2 AND 3 HERE

For comparison, Figure 3 reports results based on Diebold et al.'s (1998) test. Panel A plots the histogram of the PITs of output growth for both the density nowcast (left-hand panel) and the one-year-ahead density forecast (right-hand panel). In addition to the PITs, we also depict the 95% confidence interval (dotted lines) using a Normal approximation to a binomial distribution similar to Diebold et al.'s (1998). Both current year and one-year-ahead density forecasts of output growth in Panel A are misspecified, although misspecification is milder in the case of one-year-ahead output growth. Figure 3, Panel B, shows the histogram of the PITs for inflation. According to this test, both the density nowcast and one-year-ahead forecast overestimate tail risk. This phenomenon is more pronounced for the nowcast. Overall, the results obtained by using Diebold et al.'s (1998) test are broadly similar to those obtained by using the test that we propose in this paper, with one important exception. In the case of one-year-ahead GDP growth forecasts, our test based on Theorem 2 does not reject, whereas the Diebold et al. (1998) test does, despite the fact that both rely on assuming iid-ness of the PITs. The discrepancy in the results is most likely due to the fact that the latter test is pointwise (for each bin), whereas we jointly test the correct specification across all quantiles in the empirical distribution function: thus, in order to correctly account for the joint null hypothesis, our test has larger critical values than theirs.

Once our tests reject, it is of interest to investigate how one can improve the calibration of the density forecast. Consider, for example, SPF's predictive densities of inflation. Figure 4 plots the historical evolution of the mean probability forecasts of inflation (solid line), together with the 2.5-th, 5-th, 50-th, 95-th and 97.5-th percentiles of the predictive distribution. The picture shows that the density forecasts for both the current year and the next have evolved over time; in particular, they have become tighter towards the end of the sample. When compared to the realization (dash-dotted line), the forecast distribution looks reasonable, as the realization is contained within the 90% confidence interval throughout the sample period.

INSERT FIGURES 4 AND 5 HERE

However, the visual evidence is misleading: after studying the PITs and implementing our test (whose results are depicted in Figures 2 and 3, Panel B), we conclude that the distribution is not correctly calibrated, i.e. on average the realizations are not consistent with the ex-ante probabilities of potential outcomes. Moreover, for the case of one-year-ahead forecasts (right-hand figure in Panel B) our test finds that there are many more realizations below the mean relative to what has been anticipated. A careful observation of Figure 4 would reinforce that evidence: frequently, the realization of one-year-ahead inflation is below the forecasted mean of the distribution. Thus, the pictures suggest that the distribution is misspecified and the reason is the misspecification of the mean. To investigate this formally, one can test for forecast unbiasedness; that is, test whether $\alpha = 0$ in the regression: $y_{t+1} - \hat{y}_{t+1} = \alpha + \varepsilon_t$. The full sample estimate is $\hat{\alpha} = -0.69$ with a t-statistics of -5.38 .²⁷ The results appear to be consistent with the message in the figures.

In fact, after adjusting the mean of the distribution by adding the estimated bias (depicted in Figure 5, left panel), one obtains a well-calibrated distribution: the right panel in Figure 5 shows the results of the test for correct calibration after the (infeasible) bias adjustment, and confirms that indeed the correct specification is not rejected by our test.

To summarize, this example shows that our test can be used as a first step to determine whether the forecast density is correctly calibrated; if our test rejects the correct calibration of the forecast density, an additional analysis of the plot of the test statistic can provide guidance on the possible sources of the problem; additional (forecast rationality) tests can then verify the conjecture and help improve the calibration of the density forecast for the future (provided the source of the misspecification does not change over time).

6 Conclusions

This paper proposes new tests for predictive density evaluation. The techniques are based on Kolmogorov-Smirnov and Cramér-von Mises-type test statistics and focus both on the whole distribution as well as specific parts of it. We also propose methodologies that can be applied to multiple-step-ahead forecast horizons. Our empirical analyses uncover that both SPF output growth and inflation density forecasts are misspecified. We also investigate possible avenues that practitioners may follow in order to improve density forecasts using

²⁷The t-statistics is constructed with a Newey-West (1987) HAC estimator for the variance. If one were worried about presence of instabilities, one could alternatively apply an unbiasedness test robust to instabilities – see Rossi and Sekhposyan (2016).

our test results.

Note that our test has wide applicability: the forecast density to be evaluated in our framework can be obtained in many ways, either frequentist or Bayesian. In fact, in this paper, we are proposing a general way to evaluate forecast distributions; in particular, if one is interested in evaluating whether a forecast distribution obtained by any method (including Bayesian methods and Bayesian model averaging) is correctly specified using frequentist methods, one can use the method we propose.²⁸

References

- [1] Aastveit, K.A., C. Forni and F. Ravazzolo (2017), “Density Forecasts with MIDAS Models”, *Journal of Applied Econometrics* 32(4), 783-801.
- [2] Amisano, G. and R. Giacomini (2007), “Comparing Density Forecasts via Weighted Likelihood Ratio Tests,” *Journal of Business and Economic Statistics* 25(2), 177-190.
- [3] Ang, A., G. Bekaert and M. Wei (2007), “Do Macro Variables, Asset Markets or Surveys Forecast Inflation Better?” *Journal of Monetary Economics* 54, 1163-1212.
- [4] Bai, J. (2003), “Testing Parametric Conditional Distributions of Dynamic Models,” *Review of Economics and Statistics* 85(3), 531-549.
- [5] Bai, J. and S. Ng (2005), “Tests for Skewness, Kurtosis, and Normality for Time Series Data,” *Journal of Business and Economic Statistics* 23(10), 49-60.
- [6] Berkowitz, J. (2001), “Testing Density Forecasts, With Applications to Risk Management,” *Journal of Business and Economic Statistics* 19(4), 465-474.
- [7] Billio, M., R. Casarin, F. Ravazzolo and H. van Dijk (2013), “Time-varying Combinations of Predictive Densities using Nonlinear Filtering,” *Journal of Econometrics* 177(2), 213–232.
- [8] Bontemps, C. and N. Meddahi (2012), “Testing Distributional Assumptions: A GMM Approach,” *Journal of Applied Econometrics* 27(6), 978-1012.

²⁸For example, there are several cases in the literature where a model is estimated with Bayesian methods and yet inference is based on PITs, e.g. Clark (2011).

- [9] Clark, T. (2011), “Real-Time Density Forecasts from VARs with Stochastic Volatility,” *Journal of Business and Economic Statistics* 29(3), 327-341.
- [10] Clements, M. P. (2004), “Evaluating the Bank of England Density Forecasts of Inflation,” *The Economic Journal* 114, 844–866.
- [11] Corradi, V., S. Jin and N. Swanson (2016), “Improved Tests for Forecast Comparisons”, *mimeo*, University of Surrey.
- [12] Corradi, V. and N. R. Swanson (2006a), “Bootstrap Conditional Distribution Tests in the Presence of Dynamic Misspecification,” *Journal of Econometrics* 133, 779-806.
- [13] Corradi, V. and N. R. Swanson (2006b), “Predictive Density Evaluation,” In: G. Elliott, C. Granger and A. Timmermann (eds.), *Handbook of Economic Forecasting Vol. 1*, Elsevier, 197-284.
- [14] Corradi, V. and N. R. Swanson (2006c), “Predictive Density and Conditional Confidence Interval Accuracy Tests,” *Journal of Econometrics* 135(1–2), 187-228.
- [15] Croushore, D. and T. Stark (2001), “A Real-time Data Set for Macroeconomists,” *Journal of Econometrics* 105(1), 111-130.
- [16] Davidson, J. (1994), *Stochastic Limit Theory: An Introduction for Econometricians*, Oxford University Press.
- [17] Diks, C., V. Panchenkob and D. van Dijk (2011), “Likelihood-based Scoring Rules for Comparing Density Forecasts in Tails,” *Journal of Econometrics* 163, 215–230.
- [18] Diebold, F. X., T. A. Gunther, and A. S. Tay (1998), “Evaluating Density Forecasts with Applications to Financial Risk Management,” *International Economic Review* 39(4), 863-883.
- [19] Diebold F.X., A.S. Tay and K.F. Wallis (1999), “Evaluating Density Forecasts of Inflation: the Survey of Professional Forecasters.” In: Engle R.F. and H. White, *Cointegration, Causality, and Forecasting: A Festschrift in Honour of Clive W.J. Granger*, Oxford University Press, 76-90.
- [20] Durbin, J. (1973), *Distribution Theory for Tests Based on the Sample Distribution Function*. Philadelphia: SIAM.

- [21] Edge, R. M. and R. S. Gürkaynak (2010), “How Useful Are Estimated DSGE Model Forecasts for Central Bankers?” *Brookings Papers on Economic Activity* 41(2), 209-259.
- [22] Edge, R. M., M. T. Kiley and J. P. Laforge (2010), “A Comparison of Forecast Performance Between Federal Reserve Staff Forecasts, Simple Reduced-form Models, and a DSGE Model,” *Journal of Applied Econometrics* 25(4), 720-754.
- [23] Elliott, G. and A. Timmermann (2008), “Economic Forecasting,” *Journal of Economic Literature* 46, 3-56.
- [24] Elliott, G. and A. Timmermann (2016), *Economic Forecasting*, Princeton: Princeton University Press.
- [25] Franses, P. H. and D. van Dijk (2003), “Selecting a Nonlinear Time Series Model using Weighted Tests of Equal Forecast Accuracy,” *Oxford Bulletin of Economics and Statistics* 65, 727–744.
- [26] Giacomini, R. and H. White (2006), “Tests of Conditional Predictive Ability”, *Econometrica* 74(6), 1545-1578.
- [27] González-Rivera, G. and E. Yoldas (2012), “Autocontour-based evaluation of Multivariate Predictive Densities,” *International Journal of Forecasting* 28(2), 328-342.
- [28] González-Rivera, G. and Y. Sun (2014), “Generalized Autocontours: Evaluation of Multivariate Density Models,” *International Journal of Forecasting*, forthcoming.
- [29] Gürkaynak, R. S., B. Kisacikoglu and B. Rossi (2013), “Do DSGE Models Forecast More Accurately Out-of-Sample than VAR Models?” In: T. Fomby, L. Kilian and A. Murphy (eds.), *Advances in Econometrics: VAR Models in Macroeconomics –New Developments and Applications Vol. 31*, forthcoming.
- [30] Hayashi, F. (2000), *Econometrics*, Princeton University Press.
- [31] Hong, Y. M. and H. Li (2005), “Nonparametric Specification Testing for Continuous Time Models with Applications to Term Structure of Interest Rates,” *Review of Financial Studies* 18(1), 37-84.
- [32] Hong, Y., H. Li and F. Zhao (2007), “Can the Random Walk Model Be Beaten in Out-of-sample Density Forecasts? Evidence From Intraday Foreign Exchange Rates,” *Journal of Econometrics* 141(2), 736–776.

- [33] Inoue, A. (2001), "Testing for Distributional Change in Time Series," *Econometric Theory* 17, 156-187.
- [34] Inoue A. and B. Rossi (2012), "Out-of-sample Forecast Tests Robust to the Window Size Choice," *Journal of Business and Economics Statistics* 30(3), 432-453.
- [35] Jore, A.S., J. Mitchell and S. P. Vahey (2010), "Combining Forecast Densities from VARs with Uncertain Instabilities," *Journal of Applied Econometrics* 25(4), 621-634.
- [36] Knueppel, M. (2015), "Evaluating the Calibration of Multi-Step-Ahead Density Forecasts Using Raw Moments", *Journal of Business and Economic Statistics* 33(2), 270-281.
- [37] Newey, W.K. and K.D. West (1987), "A Simple, Positive semi-definite, Heteroskedasticity and Auto-correlation Consistent Covariance Matrix," *Econometrica* 55(3), 703-708.
- [38] Persson, J. (1974), "Comments on Estimations and Tests of EEG Amplitude Distributions," *Electroencephalography and Clinical Neurophysiology* 37, 309-313.
- [39] Ravazzolo, F. and S. P. Vahey (2014), "Forecast Densities for Economic Aggregates from Disaggregate Ensembles," *Studies of Nonlinear Dynamics and Econometrics* 18(4), 367-381.
- [40] Rossi, B. and T. Sekhposyan (2016), "Forecast Rationality Tests in the Presence of Instabilities, With Applications to Federal Reserve and Survey Forecasts", *Journal of Applied Econometrics* 31(3), 507-532.
- [41] Shorack, G. R. and J. A. Wellner (1986), *Empirical Processes with Applications to Statistics*, Wiley.
- [42] Smets, F. and R. Wouters (2007), "Shocks and Frictions in US Business Cycles: A Bayesian DSGE Approach," *American Economic Review* 97(3), 586-607.
- [43] Smirnov, N . (1948), "Table for Estimating the Goodness of Fit of Empirical Distributions," *Annals of Mathematical Statistics* 19, 279-281.
- [44] Stark, T. (2013), "SPF Panelists' Forecasting Methods: A Note on the Aggregate Results of a November 2009 Special Survey," *Real-Time Data Research Center Research Discussion Paper*.

- [45] White, H. (2001), *Asymptotic Theory for Econometricians*, Revised Edition, Academic Press.
- [46] Weiss, M. S. (1973), “Modifications of the Kolmogorov-Smirnov Statistic for Use with Correlated Data,” *Journal of the American Statistical Association* 74, 872-875.

Appendix A. Proofs

This Appendix provides proofs for Theorems 1, 2, 3, 4 and Corollaries 5 and 6. The sequence of the proofs is as follows. First, we prove part (i) and (ii) in Theorem 3, then proceed to proving Theorem 1, which follows from Theorem 3(i,ii) under Assumption 2; finally, we prove Theorems 2 and part (iii) in Theorem 3, Theorem 4 and Corollaries 5 and 6.

Proof of Theorem 3(i,ii). (i) Under Assumption 1(ii) and H_0 in eq. (3), by the proof of Lemma 1 in Bai (2003), $\{z_{t+h}\}_{t=R}^T$ is $U(0, 1)$. (ii) In what follows, we show that Assumption 1(i) satisfies the assumptions in Theorem 3.49 in White (2001) as well Assumption A in Inoue (2001). If Z_t is strong α -mixing with coefficients of size $-\delta$, $\delta > 0$, so is any measurable function of a finite number of leads and lags of Z_t (White, 2001, Theorem 3.49); in our context, $g(Z_t, \dots, Z_{t-R})$ is the cumulative distribution function and R is finite by Assumption 1(iv). Furthermore, since $g(Z_t, \dots, Z_{t-R})$ is strong mixing with coefficient $\alpha(m)$ of size $-\delta$ then $\alpha(m) = O(m^{-\delta-\epsilon})$ for some $\epsilon > 0$ (White, 2001, Definition 3.45). That is, there exists a constant $B < \infty$ such that $\frac{|\alpha(m)|}{m^{-\delta-\epsilon}} \leq B$ for every m (Davidson, 1994, p.31). Assumption A in Inoue (2001) requires that $\sum_{m=1}^{\infty} m^2 \alpha(m)^{\frac{\gamma}{4+\gamma}} < \infty$ for some $\gamma \in (0, 2)$. Note that

$$\begin{aligned} \sum_{m=1}^{\infty} m^2 \alpha(m)^{\frac{\gamma}{4+\gamma}} &= \sum_{m=1}^{\infty} m^2 \left(\frac{\alpha(m)}{m^{-\delta-\epsilon}} \right)^{\frac{\gamma}{4+\gamma}} (m^{-\delta-\epsilon})^{\frac{\gamma}{4+\gamma}} \leq \sum_{m=1}^{\infty} m^2 \left| \frac{\alpha(m)}{m^{-\delta-\epsilon}} \right|^{\frac{\gamma}{4+\gamma}} (m^{-\delta-\epsilon})^{\frac{\gamma}{4+\gamma}} \\ &\leq B^{\frac{\gamma}{4+\gamma}} \sum_{m=1}^{\infty} m^2 (m^{-\delta-\epsilon})^{\frac{\gamma}{4+\gamma}} \leq \bar{B} \sum_{m=1}^{\infty} m^{2-(\delta+\epsilon)\frac{\gamma}{4+\gamma}}, \end{aligned}$$

where $\bar{B} \equiv B^{\frac{\gamma}{4+\gamma}} < \infty$. The series $\sum_{m=1}^{\infty} m^{2-(\delta+\epsilon)\frac{\gamma}{4+\gamma}}$ is a harmonic series, convergent if $2-(\delta+\epsilon)\frac{\gamma}{4+\gamma} < -1$, i.e. if $\delta > 3\frac{4+\gamma}{\gamma}$. Thus, our Assumption 1(i) satisfies Inoue’s Assumption A. Assumption 1(ii, iii) satisfy Inoue’s Assumption B under the null. Consequently, Theorem 3 follows from Inoue (2001) by letting (in Inoue’s notation) $r = 1$. ■

Proof of Theorem 1. (i) Follows from Bai (2003, lemma 1).

(ii) The result follows from Theorem 3 noting that, from Inoue (2001 p.161, letting $r = 1$ in his notation), under *iid*, the covariance simplifies to $\sigma(r_1, r_2) = F_0(r_1, r_2) - F(r_1)F(r_2) =$

$\min(r_1, r_2) - r_1 r_2$, where the last equality follows from Shorack and Wellner (1986, p.131) and the fact that $\{z_{t+1}\}_{t=R}^T$ is *iid* Uniform(0,1). ■

Proof of Theorems 2 and 3(iii). Theorem 2 follows from Theorem 1 by the Continuous Mapping theorem. Similarly, Theorem 3 follows from Theorem 3(i,ii) by the Continuous Mapping theorem. ■

Proof of Theorem 4. (i) In what follows, we show that Assumption 1(i') satisfies the assumptions in Theorem 3.49 in White (2001) as well Assumption C in Inoue (2001). If Z_t is strong α -mixing with coefficients of size $-\delta$, $\delta > 0$, so is any measurable function of a finite number of leads and lags of Z_t , $g(Z_t, \dots, Z_{t-R})$ (White, 2001, Theorem 3.49); in our context, $g(Z_t, \dots, Z_{t-R})$ is the cumulative distribution function and R is finite by Assumption 1(iv). Furthermore, since $g(Z_t, \dots, Z_{t-R})$ is strong mixing with coefficients $\alpha(m)$ of size $-\delta$ then $\alpha(m) = O(m^{-\delta-\epsilon})$ for some $\epsilon > 0$ (White, 2001, Definition 3.45). That is, there exists a constant $B < \infty$ such that $\frac{|\alpha(m)|}{m^{-\delta-\epsilon}} \leq B$ for every m (Davidson, 1994, p.31). Assumption C in Inoue (2001) requires that $\sum_{m=1}^{\infty} m^{Q/2} \alpha(m)^{\frac{\gamma}{Q+\gamma}} < \infty$ for some $\gamma > 0$ and even integer $Q \geq 16$. Note that:

$$\begin{aligned} \sum_{m=1}^{\infty} m^{Q/2} \alpha(m)^{\frac{\gamma}{Q+\gamma}} &= \sum_{m=1}^{\infty} m^{Q/2} \left(\frac{\alpha(m)}{m^{-\delta-\epsilon}} \right)^{\frac{\gamma}{Q+\gamma}} (m^{-\delta-\epsilon})^{\frac{\gamma}{Q+\gamma}} \leq \sum_{m=1}^{\infty} m^{Q/2} \left| \frac{\alpha(m)}{m^{-\delta-\epsilon}} \right|^{\frac{\gamma}{Q+\gamma}} (m^{-\delta-\epsilon})^{\frac{\gamma}{Q+\gamma}} \\ &\leq B^{\frac{\gamma}{Q+\gamma}} \sum_{m=1}^{\infty} m^{Q/2} (m^{-\delta-\epsilon})^{\frac{\gamma}{Q+\gamma}} \leq \bar{B} \sum_{m=1}^{\infty} m^{Q/2 - (\delta+\epsilon)\frac{\gamma}{Q+\gamma}}, \end{aligned}$$

where $\bar{B} \equiv B^{\frac{\gamma}{Q+\gamma}} < \infty$. The series $\sum_{m=1}^{\infty} m^{\frac{Q}{2} - (\delta+\epsilon)\frac{\gamma}{Q+\gamma}}$ is a harmonic series, convergent if $\frac{Q}{2} - (\delta + \epsilon) \frac{\gamma}{Q+\gamma} < -1$, i.e. if $\delta > \left(1 + \frac{Q}{2}\right) \frac{Q+\gamma}{\gamma}$. Thus, our Assumption 1(i') satisfies Inoue's Assumption C. Assumption 1(ii, iii) satisfy Inoue's Assumption B under the null. Consequently, Theorem 3 follows from Inoue (2001) by letting (in Inoue's notation) $r = 1$. ■

Proof of Corollaries 5 and 6. Assumption 1(i'') is the same as assumption (i) in Theorem 3 in Giacomini and White (2006). If Z_t is strong mixing with coefficients $\alpha(m)$ of size $-\lambda/(\lambda - 2)$, so is any measurable function of a finite number of leads and lags of Z_t , $g(Z_t, \dots, Z_{t-R})$, where R is finite (White, 2001, Theorem 3.49). For our purposes, $g(Z_t, \dots, Z_{t-R}) = \mathcal{H}(z_{t+h})$. In Berkowitz's case, let $\mathcal{H}(z_{t+h}) = \oplus^{-1}(z_{t+h})$; in Knueppel's case, $\mathcal{H}(z_{t+h})$ is a general function and, for practical purposes, he recommends either $\mathcal{H}(z_{t+h}) = \oplus^{-1}(z_{t+h})$ or $\mathcal{H}(z_{t+h}) = [\sqrt{12}(z_{t+h} - \frac{1}{2})]$. Note that, invoking White (2001, Theorem 3.49), $\mathcal{H}(z_{t+h})^r$ is also mixing of size $-\lambda/(\lambda - 2)$. Then, under Assumption 4, a standard Central Limit Theorem applies as in Giacomini and White (2006, Theorem 3) and the corollaries follow directly. ■

Appendix B. Tables and Figures

Table 1. Critical Values

		κ_α			C_α			
		$\alpha :$	0.01	0.05	0.10	0.01	0.05	0.10
Panel A. Tests on the Whole Distribution								
Correct Specification Test			1.61	1.34	1.21	0.74	0.46	0.35
Panel B. Tests on Specific Parts of the Distribution								
Left Tail	$r \in [0, 0.25]$		1.24	1.00	0.88	0.56	0.34	0.24
Left Half	$r \in [0, 0.50]$		1.54	1.26	1.12	0.86	0.52	0.38
Right Half	$r \in [0.50, 1]$		1.53	1.25	1.12	0.85	0.52	0.38
Right Tail	$r \in [0.75, 1]$		1.24	1.00	0.88	0.56	0.34	0.24
Center	$r \in [0.25, 0.75]$		1.61	1.33	1.19	1.18	0.71	0.52
Tails	$r \in \{[0, 0.25] \cup [0.75, 1]\}$		1.33	1.10	0.99	0.41	0.27	0.21

Note: Panel A reports the critical values for the test statistics κ_P and C_P at the 1%, 5% and 10% nominal sizes ($\alpha = 0.01, 0.05$ and 0.10). Panel B reports the critical values for the same statistics for specific parts of the distributions, indicated in the first and second columns. The number of Monte Carlo replications is 1,000,000. The domain for r is discretized with $\underline{r} = [0 : 0.001 : 1]$.

Table 2: Size Properties

Panel A: DGP S1 (Asymptotic Critical Values)									
P	$R :$	κ_P				C_P			
		25	50	100	200	25	50	100	200
25		0.044	0.046	0.048	0.049	0.047	0.056	0.055	0.052
50		0.046	0.050	0.045	0.046	0.046	0.049	0.051	0.050
100		0.050	0.051	0.045	0.049	0.050	0.053	0.047	0.053
200		0.051	0.050	0.055	0.051	0.054	0.055	0.053	0.052
500		0.061	0.056	0.053	0.056	0.058	0.054	0.055	0.053
1000		0.053	0.050	0.048	0.047	0.050	0.051	0.050	0.050

Panel B: DGP S1 (Finite Sample Critical Values)									
P	$R :$	κ_P				C_P			
		25	50	100	200	25	50	100	200
25		0.050	0.054	0.057	0.053	0.048	0.059	0.062	0.053
50		0.050	0.059	0.048	0.050	0.044	0.057	0.052	0.052
100		0.050	0.047	0.045	0.049	0.051	0.052	0.049	0.052
200		0.049	0.048	0.057	0.051	0.052	0.053	0.059	0.054
500		0.057	0.052	0.056	0.049	0.058	0.055	0.061	0.050
1000		0.052	0.053	0.055	0.050	0.045	0.055	0.053	0.053

Note: The table reports empirical rejection frequencies for the test statistics κ_P and C_P in eqs. (4) and (5) at the 5% nominal size for various values of P and R . The number of Monte Carlo replications is 5,000. The domain for r is discretized: $\underline{r} = [0 : 0.001 : 1]$. The critical values used for Panel A are based on the asymptotic distribution, reported in Table 1, Panel A, while the critical values used for Panel B are based on simulated finite-sample distributions with 5,000 Monte Carlo replications.

Table 3: Size Properties

Panel A. DGP S2 (IID Case)									
		κ_P				C_P			
P	$R :$	25	50	100	200	25	50	100	200
25		0.044	0.046	0.044	0.042	0.049	0.048	0.049	0.050
50		0.051	0.046	0.044	0.048	0.054	0.053	0.049	0.052
100		0.046	0.044	0.044	0.046	0.051	0.047	0.046	0.047
200		0.045	0.051	0.050	0.050	0.044	0.051	0.047	0.049
500		0.050	0.048	0.050	0.053	0.048	0.046	0.048	0.051
1000		0.046	0.045	0.048	0.044	0.050	0.045	0.046	0.049
Panel B. DGP S3 (GARCH Case)									
		κ_P				C_P			
P	$R :$	25	50	100	200	25	50	100	200
25		0.046	0.046	0.048	0.049	0.049	0.056	0.054	0.052
50		0.048	0.051	0.045	0.046	0.048	0.050	0.051	0.049
100		0.051	0.052	0.045	0.050	0.052	0.054	0.047	0.054
200		0.051	0.050	0.055	0.051	0.054	0.056	0.053	0.052
500		0.061	0.057	0.053	0.056	0.057	0.056	0.055	0.054
1000		0.054	0.050	0.048	0.047	0.050	0.051	0.050	0.050
Panel C. DGP S4 (Serially Correlated Case)									
		κ_P				C_P			
P	$R :$	25	50	100	200	25	50	100	200
25		0.130	0.135	0.140	0.128	0.139	0.138	0.143	0.141
50		0.055	0.060	0.065	0.058	0.061	0.058	0.067	0.067
100		0.061	0.058	0.057	0.062	0.061	0.063	0.062	0.063
200		0.066	0.066	0.065	0.063	0.069	0.064	0.066	0.067
500		0.071	0.066	0.058	0.061	0.069	0.065	0.059	0.059
1000		0.061	0.063	0.060	0.063	0.063	0.060	0.063	0.062

Note: The table reports empirical rejection frequencies for the test statistics κ_P and C_P in eqs. (4) and (5) at the 5% nominal size. The number of Monte Carlo replications is 5,000 and $\underline{r} = [0 : 0.001 : 1]$. Critical values for DGPs S2-S3 are in Table 1, Panel A. For DGP S4 and S5, the critical values are based on 200 bootstrap replications with block length 12 in all panels, except Panel G, where it is 8.

Table 3 (continued): Size PropertiesPanel D: DGP S5 (IMA Case, $h = 2$)

		κ_P				C_P			
P	$R :$	25	50	100	200	25	50	100	200
25		0.094	0.094	0.095	0.100	0.097	0.098	0.101	0.101
50		0.043	0.041	0.044	0.041	0.038	0.037	0.042	0.038
100		0.044	0.046	0.047	0.044	0.040	0.041	0.042	0.041
200		0.052	0.053	0.053	0.054	0.043	0.049	0.044	0.045
500		0.048	0.050	0.058	0.051	0.046	0.047	0.051	0.048
1000		0.048	0.051	0.047	0.050	0.048	0.045	0.044	0.051

Panel E: DGP S5 (IMA Case, $h = 4$)

		κ_P				C_P			
P	$R :$	25	50	100	200	25	50	100	200
25		0.094	0.097	0.099	0.096	0.096	0.095	0.099	0.090
50		0.037	0.040	0.043	0.033	0.033	0.037	0.041	0.029
100		0.039	0.037	0.044	0.039	0.034	0.038	0.036	0.032
200		0.042	0.045	0.048	0.046	0.035	0.041	0.040	0.038
500		0.045	0.049	0.050	0.047	0.042	0.045	0.046	0.042
1000		0.047	0.048	0.043	0.046	0.042	0.041	0.040	0.044

Panel F: DGP S5 (IMA Case, $h = 12$)

		κ_P				C_P			
P	$R :$	25	50	100	200	25	50	100	200
25		0.091	0.095	0.093	0.093	0.089	0.095	0.086	0.098
50		0.038	0.040	0.038	0.037	0.033	0.035	0.035	0.033
100		0.038	0.044	0.045	0.038	0.030	0.033	0.033	0.031
200		0.040	0.044	0.046	0.040	0.036	0.037	0.036	0.035
500		0.044	0.045	0.053	0.047	0.036	0.044	0.047	0.039
1000		0.046	0.047	0.043	0.047	0.043	0.040	0.040	0.044

Panel G: DGP S5 (IMA Case - Robustness, $h = 12$)

		κ_P				C_P			
P	$R :$	25	50	100	200	25	50	100	200
25		0.038	0.045	0.044	0.043	0.035	0.038	0.038	0.037
50		0.037	0.038	0.037	0.040	0.028	0.034	0.029	0.028
100		0.038	0.042	0.045	0.037	0.028	0.034	0.032	0.030
200		0.038	0.040	0.043	0.038	0.034	0.034	0.031	0.034
500		0.041	0.043	0.050	0.040	0.033	0.038	0.042	0.034
1000		0.040	0.046	0.039	0.043	0.034	0.035	0.036	0.039

Table 4. Power Properties

c	DGP P1		ν	DGP P2	
	κ_P	C_P		κ_P	C_P
0	0.064	0.064	30	0.041	0.049
0.15	0.061	0.063	10	0.069	0.056
0.30	0.082	0.086	7	0.144	0.103
0.35	0.131	0.136	6	0.260	0.190
0.40	0.261	0.257	5	0.501	0.459
0.45	0.548	0.528	4	0.848	0.867
0.50	0.865	0.882	3	0.998	0.998
0.60	1.000	1.000	2	1.000	1.000

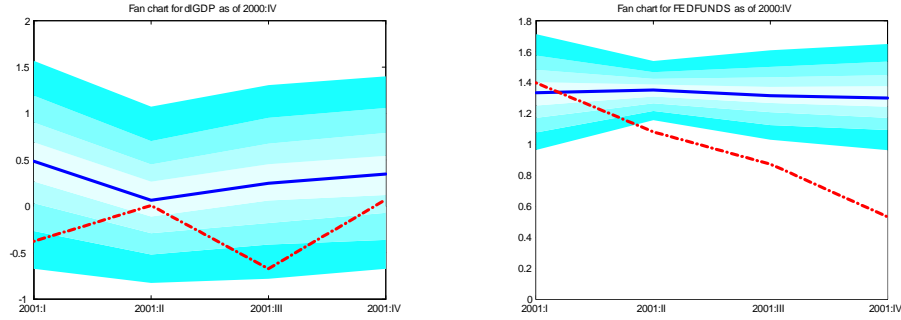
Note: The table reports empirical rejection frequencies for the test statistics κ_P and C_P in eqs. (4) and (5) for $P = 960$ and $R = 40$; the nominal size is 5%. The number of Monte Carlo replications is 5,000. The domain for r is discretized: $\underline{r} = [0 : 0.001 : 1]$. We use the critical values reported in Table 1, Panel A, to calculate the empirical rejection frequencies.

Table 5: Correct Specification Tests for SPF's Probability Forecasts

Series Name:	GDP Growth		GDP Deflator Growth	
	κ_P	C_P	κ_P	C_P
Current-year Forecasts	1.534*	0.773*	3.979*	5.066*
One-year-ahead Forecasts	0.821	0.110	3.734*	5.890*

Note: Asterisks ‘*’ indicate rejection at 5% significance level based on the critical values in Theorem 2 (reported in Table 1, Panel A). The domain for r is discretized: $\underline{r} = [0 : 0.001 : 1]$.

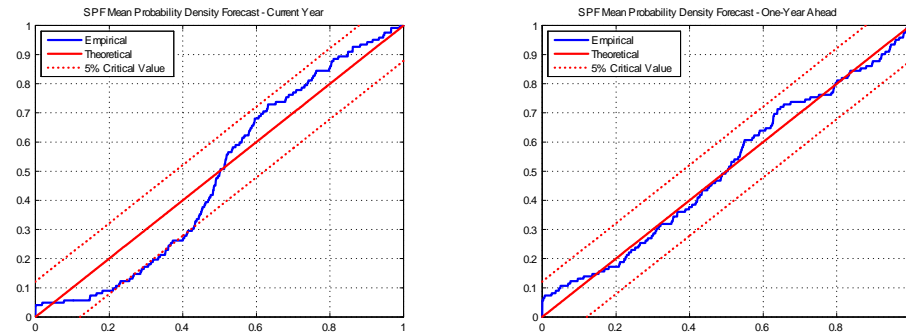
Figure 1. Fan Charts from a Representative Model in 2000:IV



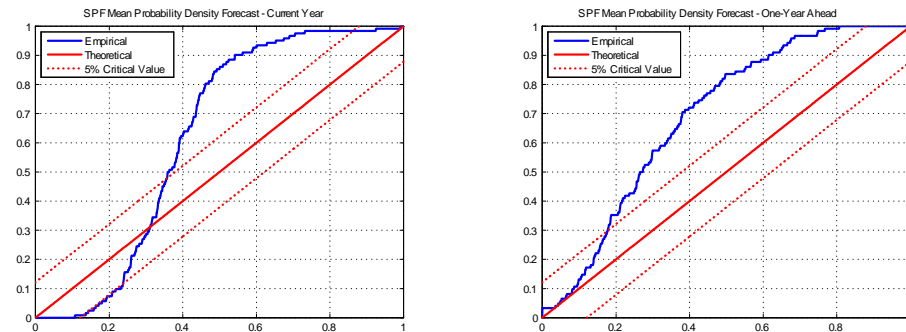
Note: The figure shows fan charts obtained by estimating the Smets and Wouters' (2007) model with data up to 2000:IV, prior to the 2001:II-2001:IV recession. The shaded areas depict the deciles of the forecast distribution for one- to four-quarter-ahead out-of-sample forecasts. The solid line represents the median forecast, while the dash-dotted line represents the actual realizations of the data.

Figure 2. CDF of the PITs – SPF Probability Forecast

Panel A: GDP Growth



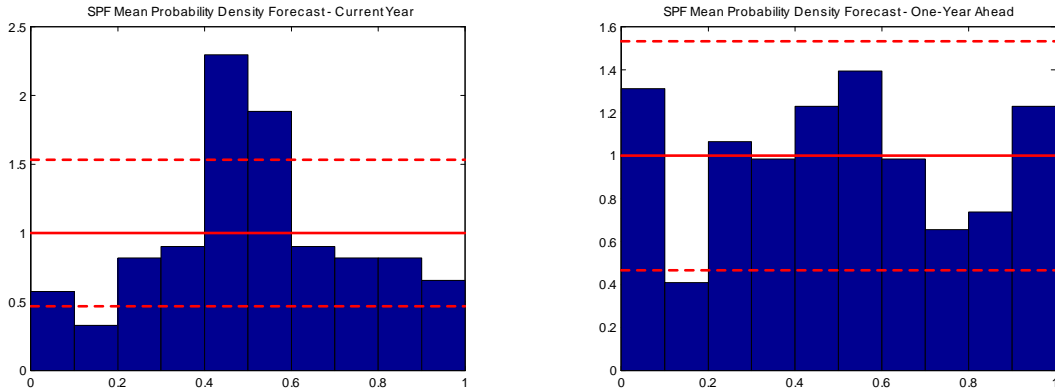
Panel B: GDP Deflator Growth



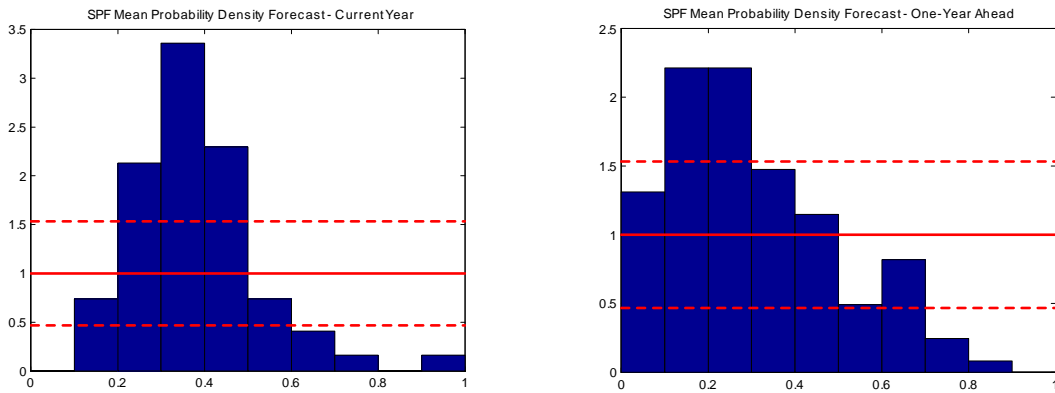
Note: The figure shows the empirical CDF of the PITs (solid line), the CDF of the PITs under the null hypothesis (the 45 degree line) and the 5% critical values bands based on the κ_P test reported in Table 1, Panel A.

Figure 3. Histogram of the PITs – SPF Probability Forecast

Panel A: GDP Growth

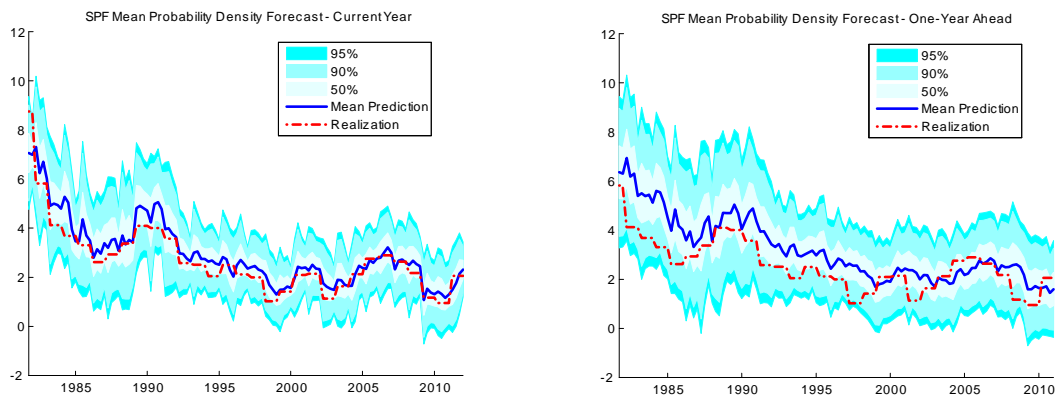


Panel B: GDP Deflator Growth



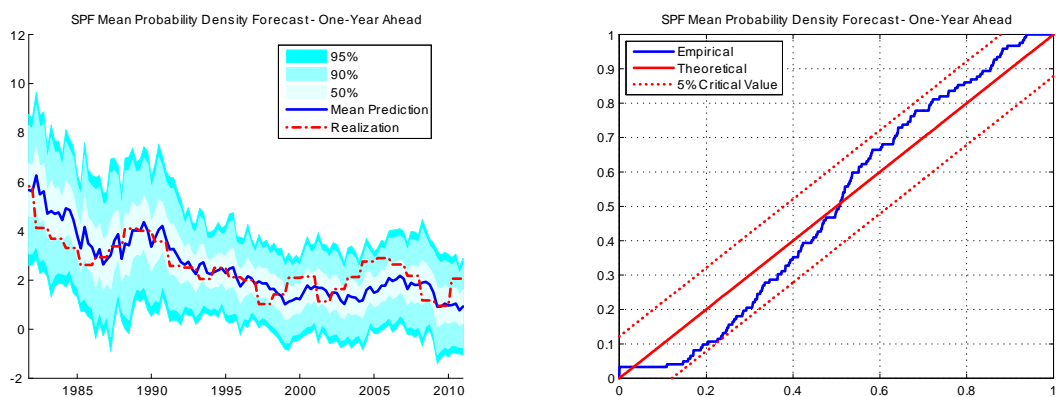
Note: The figures show the histograms of the PITs (normalized) and the 95% confidence interval approximated by Diebold et al.'s (1998) binomial distribution (dashed lines), constructed using a Normal approximation.

Figure 4. Mean and Quantiles of the SPF Inflation Density Forecast



Note: The figures plot quantiles of the SPF density forecast over time. The quantiles are constructed based on a normality assumption on the average SPF density forecasts at each point in time.

Figure 5: Forecast Evaluation of Bias-adjusted SPF Inflation



Notes: The figures show quantiles of the SPF density forecast of next year's inflation (left panel), as well as the test of correct calibration (right panel) after the correction to account for the average bias of the SPF in the observed sample.