

# Survey density forecast comparison in small samples

Laura Coroneo\*  
University of York

Fabrizio Iacone  
Università degli Studi di Milano  
University of York

Fabio Profumo  
University of York

19th June 2023

## Abstract

We apply fixed- $b$  and fixed- $m$  asymptotics to tests of equal predictive accuracy and of encompassing for survey density forecasts. We verify in an original Monte Carlo design that fixed-smoothing asymptotics delivers correctly sized tests in this framework, even when only a small number of out of sample observations is available. We use the proposed density forecast comparison tests with fixed-smoothing asymptotics to assess the predictive ability of density forecasts from the European Central Bank's Survey of Professional Forecasters (ECB SPF). We find an improvement in the predictive ability of the ECB SPF since 2010, suggesting a change in the forecasting practice after the financial crisis.

*Keywords:* survey density forecast comparison, ECB SPF, Diebold-Mariano test, forecast encompassing, fixed-smoothing asymptotics

*JEL Classification:* C12, C22, E17

---

\*Corresponding author: Laura Coroneo, Department of Economics and Related Studies, University of York, Heslington, York, YO10 5DD, UK. The support of the ESRC grant ES/J500215/1 is gratefully acknowledged. We thank Liudas Giraitis, Alastair Hall, and participants to the Sheffield Advances in Econometrics Workshop, the 42nd International Symposium on Forecasting (ISF 2022 - University of Oxford) and the 10th Italian Congress of Econometrics and Empirical Economics (ICEEE - University of Cagliari) for useful comments.

# 1 Introduction

Expectations play a key role in economic decision-making and largely determine policy outcomes. This is particularly true for monetary policy, as its effects heavily depend on expectations. For this reason, central banks around the world regularly run surveys of professional forecasters to gather information about private agents' expectations.

Survey respondents are asked to report their point forecasts for a set of macroeconomic fundamentals and, increasingly, to provide a density forecast that describes the predicted probability distribution of the variables of interest. Compared to the more popular point forecasts, density forecasts provide a wider understanding of the uncertainty associated with the prediction, see Fair (1980) and Dawid (1984) for some early references, and Tay and Wallis (2000) for a more recent detailed discussion.

Well-known examples of survey density forecasts include the Survey of Professional Forecasters (SPF) currently managed by the Federal Reserve Bank of Philadelphia, the Survey of External Forecasters managed by the Bank of England and the European Central Bank's Survey of Professional Forecasters (ECB SPF). A large amount of work has been devoted to analysing the density forecasts provided by the US SPF, see among others Diebold, Tay and Wallis (1999) and Clements (2014), and the Bank of England's Survey of External Forecasters, see among others Boero, Smith and Wallis (2008) and Mitchell and Hall (2005). The literature dedicated to density forecasts provided by the ECB SPF is more limited, see de Vincent-Humphreys, Dimitrova, Falck and Henkel (2019) for a survey, possibly because the ECB SPF started only recently, in 1999.

A challenge in forecast comparison studies for survey data is that traditional inference methods suffer from relevant small sample size distortions, which can lead to spurious results, as well documented by Clark (1999) for the Diebold and Mariano (1995) equal predictive accuracy test. This shortcoming is of course especially relevant when the analysis is performed on subsamples, as for example when only using the post-great financial crisis sample.

In this paper, we apply fixed- $b$  and fixed- $m$  asymptotics to address the small sample bias of density forecast comparison tests. We compare alternative density forecasts by testing two null hypotheses. The first hypothesis is the null of equal predictive accuracy of two forecasts, this is the Diebold and Mariano (1995) equal predictive ability test. The second one is the null of forecast encompassing in Harvey et al. (1998), which involves testing whether one forecast is encompassed by the other. To accommodate forecasts reported as probabilities for intervals, or bins, as typical for survey forecasts, we use two loss functions: the Quadratic Probability Score by Brier (1950) and the Ranked Probability Score by Epstein (1969). With these loss functions, we show that both tests can be performed in the framework of semiparametric inference on the mean of a process. In the case of the test of equal predictive accuracy, this coincides with the framework in Diebold and Mariano (1995), so we will loosely refer to it as the DM framework in the remainder of the paper, even when we apply it to the forecast encompassing test.

The DM framework is particularly appealing as it is simple and the test statistic is easy to compute. To overcome the small sample bias of the DM framework, we use an alternative approach based on fixed-smoothing asymptotics. In particular, we consider fixed- $b$  asymptotics by Kiefer and Vogelsang (2005) and fixed- $m$  asymptotics by Hualde and Iacone (2017). This approach proved capable of eliminating size distortion in the equal predictive accuracy test for comparing point forecasts, see Coroneo and Iacone (2020). In an original Monte Carlo exercise, we first document that standard asymptotics deliver unreliable density forecast comparison tests in small samples, and we then verify that fixed- $b$  and fixed- $m$  asymptotics can be used with success to perform tests of equal predictive accuracy and encompassing for density forecasts.

We apply the proposed density forecast comparison tests to assess the accuracy of the ECB SPF density forecasts for three key macroeconomic variables (real GDP growth, inflation and the unemployment rate). We are interested in establishing whether ECB SPF density forecasts can beat and/or encompass simple benchmarks, such as

an unconditional Gaussian, a Gaussian distribution based on the assumption that the target variable follows a random walk without drift, and a naive forecast taken from the previous round of ECB SPF forecasts. All benchmark forecasts are produced in real-time, by using the same information available to professional forecasters at each survey deadline.

Results indicate that ECB SPF density forecasts for unemployment and real GDP growth outperformed and sometimes encompassed the benchmarks, especially at one-year ahead and in the second subsample. On the other hand, survey forecasts for inflation do not easily outperform nor encompass the benchmarks. For all the variables, however, we find evidence of an improvement in predictive ability since 2010, supporting the anecdotal evidence of a change in the forecasting practice after the financial crisis. We also find that the ECB SPF easily outperforms and encompasses the naive benchmark, indicating that professional forecasters update their information set when making their predictions and that previous round forecasts are uninformative.

This paper contributes to the literature on forecast evaluation by introducing fixed-smoothing asymptotics to density forecast comparison tests. This type of asymptotics is becoming popular for point forecast comparison, see Choi and Kiefer (2010), Harvey et al. (2017), Li and Patton (2018), Coroneo and Iacone (2020), Coroneo et al. (2022), but, to the best of our knowledge, their properties for density forecast comparison tests have not been analysed. Our novel Monte Carlo exercise confirms the small sample bias of standard density forecast comparison tests, and indicates that fixed-smoothing asymptotics successfully addresses this issue. We also contribute to the literature on forecast encompassing by showing how the forecast encompassing test for density forecasts can be implemented in the DM framework: Clements and Harvey (2010) introduce it for dichotomic variables but we extend it to continuously distributed variables.

The remainder of the paper is organised as follows. In section 2 we describe how to perform tests of equal predictive accuracy and encompassing for survey density forecasts.

In section 3, we show how to apply fixed-smoothing asymptotics to these tests. We investigate the properties of the tests in Section 4, where we present a Monte Carlo exercise and provide recommendations for the bandwidths. In Section 5 we carry out the empirical study, and in Section 6 we conclude.

## 2 Density forecast comparison

We compare two  $h$ -step ahead density forecasts made at time  $t - h$  for the variable of interest  $y_t$  using loss functions. The  $h$ -step ahead survey density forecast  $i$  specifies the probability that the variable of interest  $y_t$  falls in bin  $k$  given the information available at time  $t - h$ ,  $\mathbf{f}_{t,i} = [f_{t,i}^1, \dots, f_{t,i}^k, \dots, f_{t,i}^K]'$ , where  $f_{t,i}^k = P_{t-h,i}(y_t \in k)$  for  $k = 1, \dots, K$ .

The vector of realisations is  $\mathbf{y}_t = [y_t^1, \dots, y_t^k, \dots, y_t^K]'$ , where the indicator variable  $y_t^k = I(y_t \in k)$  takes the value of 1 if the outcome at time  $t$  falls in bin  $k$  and zero otherwise, so that  $K - 1$  elements of  $\mathbf{y}_t$  are set to 0 and one takes value 1. The forecast error is then  $\mathbf{e}_{t,i} = \mathbf{y}_t - \mathbf{f}_{t,i}$ .

The cumulative distribution function of the density forecast is  $\mathbf{F}_{t,i} = [F_{t,i}^1, \dots, F_{t,i}^k, \dots, F_{t,i}^K]'$ , where  $F_{t,i}^k = \sum_{l=1}^k f_{t,i}^l$ , and the cumulative outcome variable is  $\mathbf{Y}_t = [Y_t^1, \dots, Y_t^k, \dots, Y_t^K]'$ , where  $Y_t^k = \sum_{l=1}^k y_t^l$ . Finally, the cumulative forecast error is given by  $\mathbf{E}_{t,i} = \mathbf{Y}_t - \mathbf{F}_{t,i}$ .

We consider two loss functions that naturally accommodate forecasts reported as histograms: the Quadratic Probability Score by Brier (1950) and the Ranked Probability Score by Epstein (1969). The Quadratic Probability Score (QPS) associated with each forecast is given by

$$QPS_{t,i} = \sum_{k=1}^K (y_t^k - f_{t,i}^k)^2 = \mathbf{e}_{t,i}' \mathbf{e}_{t,i}. \quad (1)$$

This loss function penalizes equally any probability assigned to events that do not occur. As a consequence, forecasts that assign a large probability in a neighbourhood of the realised outcome are treated in the same way as forecasts that assign a small probability to that same neighbourhood and put more probability on very distant outcomes. This

may be appropriate in some situations; in many cases, however, it is desirable to consider the forecast clustering more probability in the intervals near the realised outcome as more precise. For this reason, we also consider the Ranked Probability Score (RPS) associated with each forecast, given by

$$RPS_{t,i} = \sum_{k=1}^K (Y_t^k - F_{t,i}^k)^2 = \mathbf{E}'_{t,i} \mathbf{E}_{t,i}. \quad (2)$$

This loss function has the advantage of considering the overall tendency of the forecast probability density function, as it penalizes less severely density forecasts assigning relatively larger probabilities to outcomes that are close to the true outcome. Therefore, the RPS has the desirable property of being proper in the sense that encourages the forecasters to reveal their true beliefs, see Gneiting and Raftery (2007).

Another appealing property of the QPS and the RPS is that they are always defined, even when the realisation falls in a histogram bin to which the survey forecast has assigned a zero probability. On the contrary, the more popular logarithmic score would be undefined in this case.

We use two approaches to compare the performance of two density forecasts. The first involves testing the null hypothesis of equal predictive accuracy of the two forecasts according to the QPS or the RPS loss function. This can be implemented with the test for equal predictive accuracy proposed by Diebold and Mariano (1995). The second approach involves testing for whether one density forecast is encompassed by the other one, in the sense that the predictive accuracy (according to the QPS or the RPS loss function) of the encompassing density forecast cannot be improved by a linear combination with the encompassed forecast. This is the forecast encompassing test and, for point forecasts, Harvey et al. (1998) show that, by redefining the loss differential, it is possible to implement it using the DM framework.

In what follows, we first show how the DM framework can be used also to perform the forecast encompassing test for density forecasts. We then discuss the limitations of

standard asymptotics when applied to the DM framework, and apply fixed-smoothing asymptotics to the DM framework for density forecast comparison.

## 2.1 Equal predictive accuracy

A test of equal predictive accuracy allows testing the null hypothesis that two alternative forecasts have equal forecasting accuracy according to a user-chosen loss function, which in the case of density forecasts can be the QPS or the RPS loss function.

Denote by  $L^i$  the loss function for  $i = 1, 2$ , so that  $L_t^i = QPS_{t,i}$  if the QPS loss is used or  $L_t^i = RPS_{t,i}$  if the RPS loss is used, and the loss differential by

$$d_t = L_t^1 - L_t^2, \quad (3)$$

the null hypothesis of equal forecasting ability is

$$H_0 : \{E(d_t) = 0\}. \quad (4)$$

## 2.2 Forecast encompassing

A forecast encompassing test involves testing whether one set of forecasts encompasses another one, in the sense that the accuracy of one set of (encompassing) forecasts  $\mathbf{f}_{t,1}$  cannot be improved through a linear combination with a second set of (encompassed) forecasts  $\mathbf{f}_{t,2}$ . To this end, we consider the density forecast combination

$$\mathbf{f}_{t,c}(\lambda) = (1 - \lambda)\mathbf{f}_{t,1} + \lambda\mathbf{f}_{t,2} \quad (5)$$

where  $\lambda$  ( $0 \leq \lambda \leq 1$ ) is a scalar and denotes the weight associated with forecasts  $\mathbf{f}_{t,2}$ . In this context,  $\mathbf{f}_{t,1}$  encompasses  $\mathbf{f}_{t,2}$  if the optimal weight in the QPS (or RPS) sense is equal to zero.

In Appendix A, we show that if we define  $d_t$  as

$$d_t = \begin{cases} \mathbf{e}'_{t,1}(\mathbf{e}_{t,1} - \mathbf{e}_{t,2}), & \text{for } QPS, \\ \mathbf{E}'_{t,1}(\mathbf{E}_{t,1} - \mathbf{E}_{t,2}), & \text{for } RPS, \end{cases} \quad (6)$$

then the null of density forecast encompassing can be expressed as

$$H_0 : \{E(d_t) = 0\}$$

and the density forecast encompassing test can be conducted against the one-sided alternative  $E(d_t) > 0$  (i.e.,  $\lambda > 0$ ), given the assumption of a non-negative combination weight.

### 2.3 Diebold-Mariano framework

In sections 2.1-2.2, we showed how both the equal predictive accuracy and the forecast encompassing tests can be performed in the framework of inference on the mean of the process  $d_t$  also in the context of density forecast evaluation. The difference between the two tests lies in how the process  $d_t$  is defined. For the test for equal predictive accuracy,  $d_t$  is defined as in (3), while for the test for density forecast encompassing  $d_t$  is defined as in (6).

Denoting the sample average as  $\bar{d} = \frac{1}{T} \sum_{t=1}^T d_t$  and the long run variance as  $\sigma_T^2 = \text{var}(\sqrt{T} \bar{d})$ , then the test statistic is  $\sqrt{T} \bar{d} / \sigma_T$ . Under regularity conditions, including mixing at a sufficient rate for  $d_t$ , sufficient moments for  $|d_t|$  and  $\sigma_T > 0$ , the assumptions as in Giacomini and White (2006) hold, yielding a central limit theorem for  $\sqrt{T} \bar{d}$ . Then, under  $H_0$ ,

$$\sqrt{T} \frac{\bar{d}}{\sigma_T} \rightarrow_d N(0, 1). \quad (7)$$

The test statistic in (7) is unfeasible as  $\sigma_T$  is unknown, but this may be replaced by



an estimate, say  $\hat{\sigma}$ . If the latter is consistent, in the sense  $\hat{\sigma} - \sigma_T = o_p(1)$ , the feasible statistic obtained in this way retains the standard normal limiting distribution.

One estimate that, under regularity conditions, fits this purpose, is the Weighted Covariance Estimate

$$\hat{\sigma}_{WCE}^2 = \hat{\gamma}_0 + 2 \sum_{j=1}^{T-1} k(j/M) \hat{\gamma}_j$$

where  $\hat{\gamma}_j = \frac{1}{T} \sum_{t=1}^{T-j} (d_t - \bar{d})(d_{t+j} - \bar{d})$  is the sample autocovariance,  $k(\cdot)$  is a kernel function and  $M$  is a bandwidth parameter. A popular kernel is the triangular (Bartlett) kernel

$$k^B(j/M) = 1 - \frac{j}{M} \quad \text{for } j \leq M$$

yielding the estimate

$$\hat{\sigma}_{WCE-B}^2 = \hat{\gamma}_0 + 2 \sum_{j=1}^M \left( \frac{M-j}{M} \right) \hat{\gamma}_j.$$

Regularity conditions to ensure consistency include  $M \rightarrow \infty$  and  $M/T \rightarrow 0$  as  $T \rightarrow \infty$ .

A second class of estimates of the long run variance is the Weighted Periodogram Estimate

$$\hat{\sigma}_{WPE}^2 = 2\pi \sum_{j=1}^{T/2} K_M(\lambda_j) I(\lambda_j) \quad (8)$$

where  $K_M(\lambda_j)$  is a symmetric kernel function, and  $I(\lambda_j) = \left| \frac{1}{\sqrt{2\pi T}} \sum_{t=1}^T d_t e^{i\lambda_j t} \right|^2$  is the periodogram of  $d_t$  computed at the Fourier frequencies  $\lambda_j = \frac{2\pi j}{T}$  for  $j = 1, \dots, T/2$ . A popular kernel in this case is the Daniell kernel

$$K_M^D(j) = \frac{1}{m} \quad \text{for } j \leq m$$

where  $m$  is a user-chosen parameter that is linked to the bandwidth  $M$  (and it is, with a slight abuse of notation, referred to as bandwidth too). This kernel is often a convenient choice, as the Daniell kernel estimate of the long run variance has a very simple formula

in the frequency domain,

$$\hat{\sigma}_{WPE-D}^2 = 2\pi \frac{1}{m} \sum_{j=1}^m I(\lambda_j). \quad (9)$$

When  $m \rightarrow \infty$  and  $m/T \rightarrow 0$  as  $T \rightarrow \infty$ , the estimate is consistent. An extension of this class of estimates is introduced in Phillips (2005), that shows that  $\sigma_T$  can be consistently estimated by regressing the series of interest on an orthonormal series (although Phillips (2005) actually only considers a constant long term variance, his argument also applies to the more general context we consider here). This orthonormal series may be a set of trigonometric polynomials, but this does not necessarily have to be the case.

Unfortunately, the DM framework is subject to severe size distortion in small and medium-sized samples, as documented, for example, in Clark (1999). Obviously, finite sample size distortion is not a problem affecting only the DM framework, it is common to any test that makes inference on the mean (or on a regression parameter) using a heteroskedasticity autocorrelation consistent estimate of the long run variance and maintaining the limit normality assumption for the standardised statistic, see for example Newey and West (1994). In fact, in any finite sample, the ratio  $M/T$  is still non-zero, and in a moderate size sample this ratio may be non-negligible. Thus, this size distortion may be more severe in the context of forecast comparisons, as in many cases the sample size is relatively small, when compared to other macro and financial applications.

### 3 Fixed-smoothing asymptotics

Neave (1970) shows that treating the ratio  $M/T$  as constant can provide a better measure of the variance of the weighted covariance estimate of a spectral estimate. Kiefer and Vogelsang (2002a,b, 2005) apply the same intuition to the problem of testing hypothesis about the mean for a weakly dependent process, deriving the distribution of the feasible test statistic when  $M/T \rightarrow b \in (0, 1]$  as  $T \rightarrow \infty$ . Under this assumption  $\hat{\sigma}^2$  is not consistent, and the test statistic has a non-standard limit distribution that depends both

on  $b$  and on the kernel choice. Because of the dependence on  $b$  of the limit distribution, this approach is often referred to as “fixed- $b$ ”.

In the context of the DM framework, for the Bartlett kernel the results of Kiefer and Vogelsang (2005) imply that, under  $H_0$  and regularity conditions, when  $M/T \rightarrow b \in (0, 1]$  as  $T \rightarrow \infty$ ,

$$\sqrt{T} \frac{\bar{d}}{\widehat{\sigma}_{WCE-B}} \rightarrow_d \Phi^B(b) \quad (10)$$

$\Phi^B(b)$  is characterised in Kiefer and Vogelsang (2005) and a cubic equation is provided for critical values.

In the frequency domain, fixed- $b$  corresponds to keeping  $m$  constant when the Daniell kernel is used. This naturally leads to considering asymptotics for fixed  $m$ . Under  $H_0$  and regularity conditions, Hualde and Iacone (2017) consider  $m$  constant as  $T \rightarrow \infty$ , in this case we have

$$\sqrt{T} \frac{\bar{d}}{\widehat{\sigma}_{WPE-D}} \rightarrow_d t_{2m}. \quad (11)$$

Sun (2013) shows that a limit of this kind also holds for the general orthonormal series variance estimator.

Fixed- $b$  and fixed- $m$  asymptotics can be heuristically understood as undersmoothing in the context of estimating the spectral density at frequency zero. For this reason, many references, for example Sun (2013), refers to them collectively as fixed-smoothing.

Monte Carlo simulations in Kiefer and Vogelsang (2005) suggest that critical values obtained using fixed- $b$  asymptotics result in better empirical size for tests. This was later justified theoretically by Sun (2014), that shows that fixed- $b$  asymptotics provides a higher order refinement. Moreover, fixed-smoothing asymptotics gives a justification (and suitable critical values) even for bandwidths that researchers would not consider when using standard asymptotics: it is even possible to choose  $M = T$  when using the weighted covariance Bartlett estimate, or to choose  $m = 1$  when using the weighted periodogram Daniell estimate. This allows a further correction in the empirical size, as in Monte Carlo simulations larger bandwidths  $M$  (smaller  $m$ ) are associated to better

empirical size. For example, Monte Carlo simulations in Coroneo and Iacone (2020) indicate that it is possible to completely eliminate the size distortion documented by Clark (1999).

**Assumption 1** *Partial sums of  $d_t$  are such that the functional central limit theorem (FCLT) holds*

$$\frac{\sqrt{T}}{T} \frac{1}{\sigma_T} \sum_{t=1}^{\lfloor rT \rfloor} d_t \Rightarrow W(r)$$

where  $\lfloor \cdot \rfloor$  denotes the integer part of a number,  $r \in [0, 1]$  and  $W(r)$  is a standard Brownian motion.

Assumption 1 is sufficient to establish the fixed-smoothing limits (10) and (11). This assumption is not primitive, but it is convenient because it may be established under a range of conditions. For example, Phillips and Solo (1992) consider linear processes of independent, identically distributed innovations of martingale difference sequences. On the other hand, Wooldridge and White (1988) consider mixing processes, thus allowing for forms of heteroskedasticity that may also induce non-stationarity, under the additional assumption that  $Var \left( \frac{\sqrt{T}}{T} \frac{1}{\sigma_T} \sum_{t=1}^{\lfloor rT \rfloor} d_t \right) \rightarrow r$ . In view of the non-linearity in the loss function, establishing a linear representation for  $d_t$  from primitive assumptions on  $f_t^{k,i}$  and  $y_t^k$  may be very challenging, whereas establishing mixing properties may be easier, especially when  $f_t^{k,i}$  and  $y_t^k$  are limited to being  $M$ -dependent processes. However, as the two classes may overlap but are not included in each other, see discussion in Phillips and Solo (1992) and Andrews (1984), we prefer the more general assumption given here, that encompasses them both.

## 4 Monte Carlo study of size and power

We analyse the empirical size and power of the tests of equal predictive accuracy and of encompassing for density forecast by means of a Monte Carlo experiment. Since Kiefer

and Vogelsang (2005), simulation studies have by now covered a fairly wide range of situations, including inference in regression models, in non-linear models, and others. We refer to Lazarus, Lewis, Stock and Watson (2018) for a recent, comprehensive study. In point forecasting, studies include Coroneo and Iacone (2020) on forecast evaluation in small samples, Harvey, Leybourne and Whitehouse (2017) on forecast encompassing, and Li and Patton (2018) on forecast evaluation in large samples.

We already noticed that simulation studies find that fixed-smoothing asymptotics yield better approximation of the empirical size, and that this improvement is stronger the larger is the bandwidth  $M$  (the smaller is  $m$ ). These works also find that the finite sample power is decreasing with the bandwidth, therefore documenting the existence of a trade-off between correct size and power. Lazarus, Lewis, Stock and Watson (2018), drawing on their extensive simulation study, recommend  $M = \lfloor 1.3T^{1/2} \rfloor$  and  $m = \lfloor 0.2T^{2/3} \rfloor$ .

In this section, we check whether the size improvements for the equal predictive ability and forecast encompassing tests still hold in the case of density comparisons. We use a rather small sample that replicates the dimension of the sample of our dataset. We also examine the issue of bandwidth selection, and compare our results with Lazarus, Lewis, Stock and Watson (2018) and Coroneo and Iacone (2020).

In our Monte Carlo study, for simplicity, we only consider the QPS loss function. We consider a sample of  $T$  observations, and we assume that the probability that the variable of interest  $y_t$  falls in bin  $k$ , for  $k = 1, 2, 3$ , is given by  $\mathbf{y}'_t = (0, 1, 0)$ . We also assume that we have two density forecasts that assign the probability that  $y_t$  falls in bin  $k$  as follows

$$\mathbf{f}'_{1,t} = (A_t, 1 - A_t, 0);$$

$$\mathbf{f}'_{2,t} = (0, 1 - B_t, B_t);$$

where

$$A_t = a_t + a_{t-1} + \dots + a_{t-Q};$$

$$B_t = b_t + b_{t-1} + \dots + b_{t-Q};$$

and  $a_t, \dots, a_{t-Q}$  are realisations from a uniform distribution in  $[0, \alpha/(Q+1)]$ ,  $b_t, \dots, b_{t-Q}$  are realisations from a uniform distribution in  $[0, \beta/(Q+1)]$ , and  $a_t, \dots, a_{t-Q}, b_t, \dots, b_{t-Q}$  are all independently distributed.

The forecast errors are then given by

$$\mathbf{e}'_{1,t} = (-A_t, A_t, 0);$$

$$\mathbf{e}'_{2,t} = (0, B_t, -B_t).$$

In this setting,  $E(\mathbf{e}'_{1,t}\mathbf{e}_{1,t}) = E(2A_t^2)$ ,  $E(\mathbf{e}'_{2,t}\mathbf{e}_{2,t}) = E(2B_t^2)$ . This means that the null hypothesis of the equal predictive ability test,  $E(d_t) = E(\mathbf{e}'_{1,t}\mathbf{e}_{1,t}) - E(\mathbf{e}'_{2,t}\mathbf{e}_{2,t}) = 0$ , follows from setting  $\alpha = \beta$ . For the forecast encompassing test,  $E(\mathbf{e}'_{1,t}(\mathbf{e}_{1,t} - \mathbf{e}_{2,t})) = 2E(A_t^2) - E(A_t)E(B_t)$ , so to obtain the null hypothesis  $E(d_t) = E(\mathbf{e}'_{1,t}(\mathbf{e}_{1,t} - \mathbf{e}_{2,t})) = 0$ , we set  $\beta = 8\frac{4+3Q}{12(Q+1)}\alpha$ . We can investigate the power of the equal predictive ability test setting  $\beta = \sqrt{\alpha^2 - 3/2 \times c/\sqrt{T}}$  as we increase the value of  $c$ .

In our experiment we set  $\alpha = \beta = 1$  for the equal predictive accuracy test and  $\alpha = 3/8$  for the forecast encompassing test, and  $Q$  up to 6, with sample size set at  $T = 40, 80$ , and we repeat the experiment for 10,000 replications. Our sample size is much smaller than the sample size of Lazarus, Lewis, Stock and Watson (2018), and it matches the dimension of the sample available for our empirical study. Indeed, checking the empirical performance in such small samples is one reason of interest in this experiment.

In Tables 1 and 2, we report the empirical size of the test with one-sided alternative  $H_1 : \{E(d_t) > 0\}$  when 5% critical values from both standard asymptotics and fixed-smoothing asymptotics are used. In columns WCE, the long run variance estimate is

Table 1: Empirical size of the test of equal predictive ability

<b>Panel A: standard asymptotics</b>								
$T$	$Q$	<b>WCE</b>			<b>WPE</b>			
		$[T^{1/3}]$	$[T^{1/2}]$	$T$	$[T^{1/4}]$	$[T^{1/3}]$	$[T^{1/2}]$	$[T^{2/3}]$
40	0	0.067	0.079	0.204	0.089	0.075	0.064	0.061
	2	0.102	0.098	0.218	0.084	0.078	0.082	0.117
	4	0.159	0.129	0.234	0.101	0.096	0.130	0.192
	6	0.197	0.152	0.244	0.109	0.117	0.176	0.232
80	0	0.060	0.068	0.201	0.085	0.069	0.058	0.056
	2	0.082	0.080	0.206	0.086	0.067	0.066	0.092
	4	0.117	0.098	0.220	0.090	0.077	0.087	0.160
	6	0.148	0.111	0.221	0.092	0.083	0.112	0.194

<b>Panel B: fixed-smoothing asymptotics</b>								
$T$	$Q$	<b>WCE</b>			<b>WPE</b>			
		$[T^{1/3}]$	$[T^{1/2}]$	$T$	$[T^{1/4}]$	$[T^{1/3}]$	$[T^{1/2}]$	$[T^{2/3}]$
40	0	0.047	0.045	0.049	0.052	0.051	0.051	0.055
	2	0.076	0.058	0.057	0.047	0.051	0.067	0.108
	4	0.129	0.083	0.070	0.058	0.069	0.113	0.183
	6	0.166	0.105	0.079	0.066	0.082	0.157	0.222
80	0	0.048	0.047	0.048	0.050	0.051	0.050	0.053
	2	0.067	0.053	0.052	0.050	0.049	0.057	0.085
	4	0.100	0.072	0.059	0.054	0.058	0.077	0.152
	6	0.131	0.082	0.061	0.053	0.062	0.100	0.188

Note: empirical size of the equal predictive ability test with standard asymptotics (panel A) and fixed-smoothing asymptotics (panel B). The theoretical size is 5%, for a one-sided alternative hypothesis.  $Q$  indicates the dependence in the process. WCE refers to the test statistic with Weighted Covariance Estimate with Bartlett kernel for the long run variance; WPE refers to the test statistic with Weighted Periodogram Estimate with Daniell kernel for the long run variance.

Table 2: Empirical size of the forecast encompassing test

<b>Panel A: standard asymptotics</b>								
$T$	$Q$	<b>WCE</b>			<b>WPE</b>			
		$\lfloor T^{1/3} \rfloor$	$\lfloor T^{1/2} \rfloor$	$T$	$\lfloor T^{1/4} \rfloor$	$\lfloor T^{1/3} \rfloor$	$\lfloor T^{1/2} \rfloor$	$\lfloor T^{2/3} \rfloor$
40	0	0.046	0.061	0.189	0.073	0.059	0.045	0.040
	2	0.075	0.067	0.203	0.062	0.050	0.052	0.094
	4	0.122	0.092	0.222	0.075	0.064	0.091	0.167
	6	0.163	0.128	0.208	0.088	0.096	0.151	0.194
80	0	0.047	0.056	0.187	0.075	0.055	0.046	0.043
	2	0.049	0.042	0.170	0.052	0.040	0.035	0.055
	4	0.088	0.068	0.182	0.066	0.050	0.055	0.120
	6	0.110	0.071	0.182	0.062	0.049	0.074	0.155

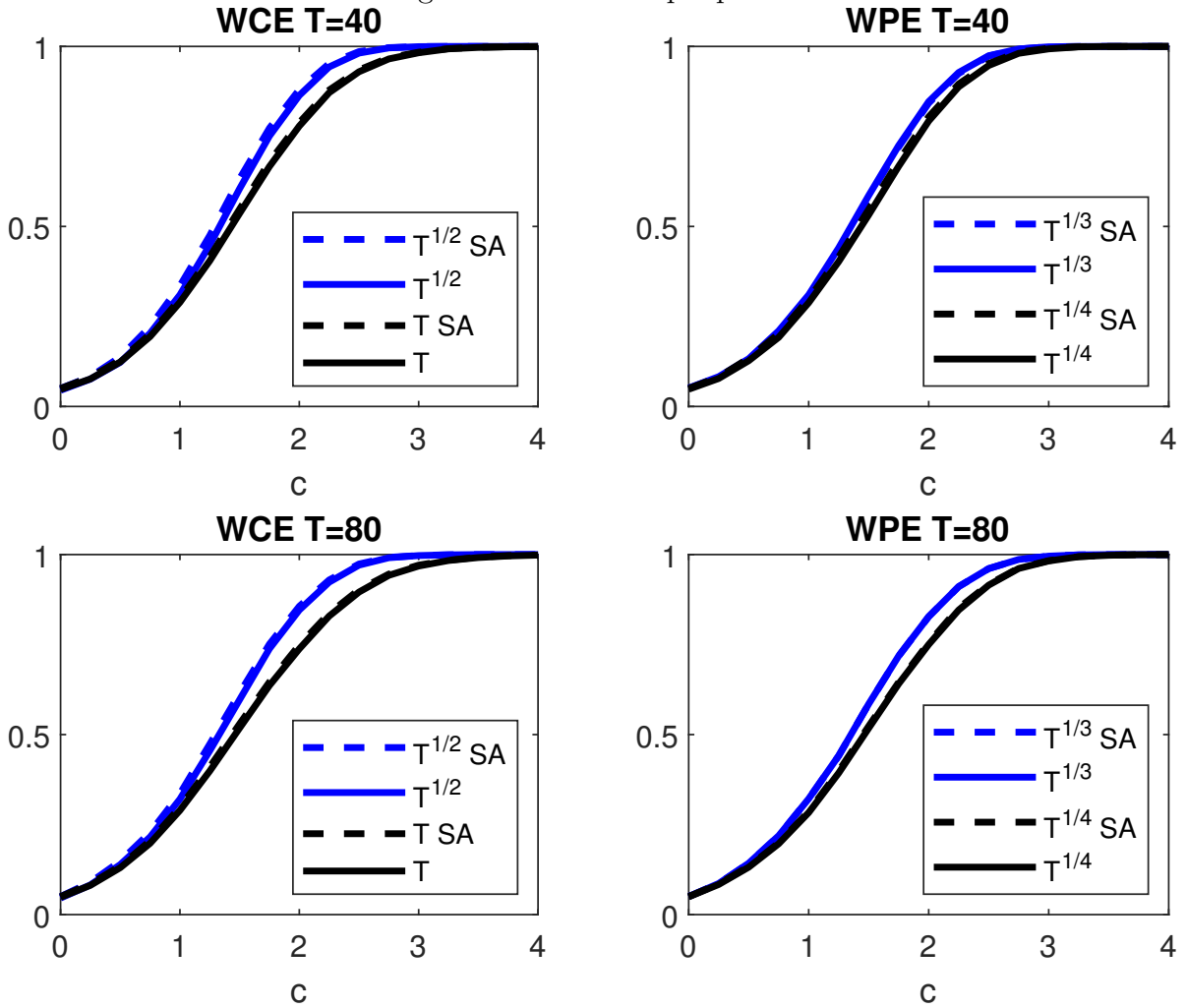
  

<b>Panel B: fixed-smoothing asymptotics</b>								
$T$	$Q$	<b>WCE</b>			<b>WPE</b>			
		$\lfloor T^{1/3} \rfloor$	$\lfloor T^{1/2} \rfloor$	$T$	$\lfloor T^{1/4} \rfloor$	$\lfloor T^{1/3} \rfloor$	$\lfloor T^{1/2} \rfloor$	$\lfloor T^{2/3} \rfloor$
40	0	0.030	0.029	0.035	0.039	0.037	0.035	0.033
	2	0.050	0.037	0.037	0.032	0.031	0.042	0.077
	4	0.092	0.051	0.057	0.038	0.042	0.077	0.157
	6	0.139	0.082	0.065	0.047	0.060	0.132	0.185
80	0	0.034	0.034	0.040	0.040	0.038	0.037	0.039
	2	0.039	0.028	0.029	0.030	0.029	0.028	0.054
	4	0.070	0.042	0.036	0.037	0.040	0.045	0.115
	6	0.090	0.050	0.040	0.035	0.035	0.062	0.148

Note: empirical size of the forecast encompassing test with standard asymptotics (panel A) and fixed-smoothing asymptotics (panel B). The theoretical size is 5%, for a one-sided alternative hypothesis.  $Q$  indicates the dependence in the process. WCE refers to the test statistic with Weighted Covariance Estimate with Bartlett kernel for the long run variance; WPE refers to the test statistic with Weighted Periodogram Estimate with Daniell kernel for the long run variance.



Figure 1: Finite sample power



Note: power performances of the equal predictive ability test in a samples of size  $T = 40$  and  $T = 80$ , for the theoretical size 5% and for a one-sided alternative hypothesis. The dashed lines refer to power performances using size-adjusted critical values while solid lines use fixed-smoothing asymptotics. The parameter  $c$  indicates the distance from the null hypothesis. WCE refers to the test statistic with Weighted Covariance Estimate with Bartlett kernel for the long run variance; WPE refers to the test statistic with Weighted Periodogram Estimate with Daniell kernel for the long run variance.

computed using a Bartlett kernel with bandwidths  $M = \lfloor T^{1/3} \rfloor$ ,  $M = \lfloor T^{1/2} \rfloor$  and  $M = T$ . In columns WPE, we use the Daniell kernel with bandwidths  $m = \lfloor T^{1/4} \rfloor$ ,  $m = \lfloor T^{1/3} \rfloor$ ,  $m = \lfloor T^{1/2} \rfloor$  and  $m = \lfloor T^{2/3} \rfloor$ . Consistently with results from other simulation studies, standard asymptotics are associated with size distortions. The performance deteriorates as the dependence increases with  $Q$ , especially when the bandwidth  $m$  is too long (or, to a lesser extent, when  $M$  is too short), reflecting the fact that the dependence causes a curvature in the spectral density at larger frequencies, and thus a bias in the estimation of

the spectral density in zero. The second source of distortion is due to the approximation of the average periodogram as its probability limit, and this is more evident when  $m$  is too short ( $m = \lfloor T^{1/4} \rfloor$ ), and when the bandwidth  $M$  is too long ( $M = T$ ). Using fixed-smoothing asymptotics always improves the empirical size. As usual, the best performance is for  $M = T$  or the smallest  $m$  (as the size distortion due to the curvature of the spectral density is least, in this case), but on balance we observe correctly sized tests with WCE already with bandwidth  $M = \lfloor T^{1/2} \rfloor$ , likewise, we observe correct size with WPE already with bandwidth  $m = \lfloor T^{1/3} \rfloor$ .

For the power study we only consider the equal predictive ability test. We set  $Q = 0$  and increasing values of  $c$  up to 4. In this case, we only consider bandwidths that are associated to good empirical size properties, namely WCE with  $M = \lfloor T^{1/2} \rfloor$  and  $M = T$ , and WPE with  $m = \lfloor T^{1/4} \rfloor$  and  $m = \lfloor T^{1/3} \rfloor$ , in all cases only for fixed-smoothing asymptotics. For the purpose of comparison only, we also plot the size adjusted power. Power performances are reported in Figure 1. In all cases the empirical power is a very good approximation of the size adjusted power, again offering support to the assumption that fixed-smoothing asymptotic is a valuable instrument for inference. We also find that, as a general rule, larger bandwidths  $M$  (smaller  $m$ ) are associated to lower power, consistently with other similar simulation studies. Overall we suggest  $M = \lfloor T^{1/2} \rfloor$  and  $m = \lfloor T^{1/3} \rfloor$ . Given our sample size, these bandwidth rules seem in line with the recommendation in Lazarus, Lewis, Stock and Watson (2018).

## 5 Application

We use the proposed density forecast evaluation tests with fixed-smoothing asymptotics to evaluate the predictive ability of density forecasts from the European Central Bank's Survey of Professional Forecasters (ECB SPF) for HICP inflation, the unemployment rate and real GDP growth against three simple benchmarks: an unconditional Gaussian density, a Gaussian random walk density, and naive benchmark that uses the latest

survey value for the same forecast horizon. Of course, one can use more sophisticated benchmarks, but here our objective is to assess whether the ECB SPF survey forecasts can at least beat three very simple benchmarks.

## 5.1 The ECB Survey of Professional Forecasters

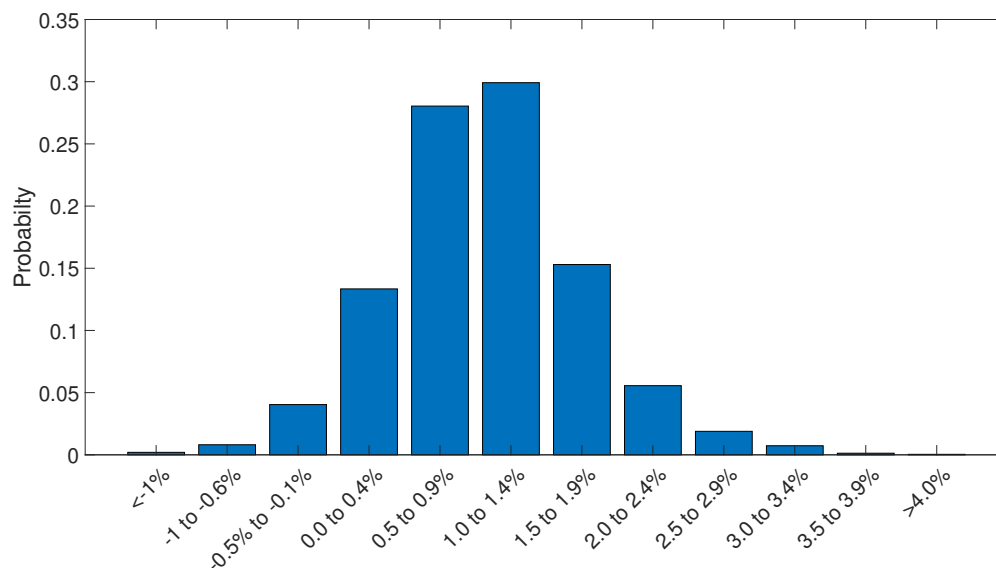
We use aggregate ECB SPF density forecasts at 1 and 2 years ahead for inflation (year-on-year percentage change of the Harmonised Index of Consumer Prices, HICP), real GDP growth (year-on-year percentage change of real GDP) and the unemployment rate (as percentage of the labour force).

The ECB SPF is administered quarterly to a panel of forecasters (about 80 institutions with an average of 60 responses each round). Participants are experts affiliated with financial or non-financial institutions based within the European Union, and form an heterogeneous group to guarantee the representativeness and independence of the expectations collected.

Participants are asked to provide a forecast for the current calendar year, the following calendar year, the calendar year after that, a long term horizon, a rolling horizon one year ahead of the latest available data and a rolling horizon two years ahead of the latest available data. For more information on the ECB SPF see Garcia (2003) and Bowles, Friz, Genre, Kenny, Meyler and Rautanen (2007).

To report their density forecasts, participants are given a set of specific ranges and are asked to predict the probability that the target variable will fall in each specific range, or bin, with the first and the last being open intervals. The number of ranges given in every survey round can change but their width is fixed. The ECB SPF reports both the anonymised individual density forecasts and the aggregate density forecast, constructed by summing up the individual probabilities reported in the SPF and dividing by the number of respondents. For example, in Figure 2 we present the one year-ahead density forecast for the December 2016 HICP produced in the 2016.Q1 survey round.

Figure 2: ECB SPF density forecast for HICP one-year ahead, December 2016



Note: histogram of the one-year ahead aggregate density forecast for HICP from the 2016.Q1 survey round. Participants are asked to report a probability for the realisation in December 2016 to fall in each bin.

## 5.2 Benchmark Forecasts

We compare ECB SPF density forecasts against three simple benchmarks: an unconditional Gaussian density forecast, a Gaussian random walk density forecast that represents a standard benchmark for forecasting, and a naive forecast based on the lagged ECB SPF density forecast that, as such, incorporates all the information available at the previous survey round.

As forecasters operate using data as available at the time the forecasts are made, we construct the Gaussian benchmark density forecasts using only the real-time information available to professional forecasters up to the deadline for responding to each survey round by using the historical vintages from the Real-time Database for the euro area built by Giannone, Henry, Lalik and Modugno (2012) and available on the European Central Bank Statistical Data Warehouse.

For the unconditional Gaussian benchmark, we use a Gaussian distribution with mean and variance obtained from the historical observations of the target variable as

available at each survey deadline. For the Gaussian random walk benchmark, we use a normal distribution with conditional expectation  $\mu_t = y_{t-h}$  and variance calculated using all historical observations as available at each survey deadline. From these normal distributions, we compute the probability that the realization of the target variable falls inside each bin.

For the naive benchmark, we simply use the last available ECB SPF density forecast for the same horizon, i.e.  $f_t^{k,Naive} = f_{t-1}^{k,SPF}$ . In the case of different bins available from a survey round to the following, the forecasts are adjusted to accommodate the new bin structure. If in the new survey round there are more bins than in the previous, the probability of the last bin is equally split across the additional bins available; if there are less bins in the current survey round than the previous round, the probabilities of extreme bins are added up and placed in the only available bin. For additional discussion about the changing bin structure see D’Amico et al. (2008) and Manzan (2021).

### 5.3 Empirical Results

We analyse the ECB SPF aggregate density forecasts at the rolling horizons of one and two years for HICP inflation, the unemployment rate and real GDP growth, for the surveys between 2000.Q1 and 2019.Q4, corresponding to a total of 80 quarterly observations. We also split the sample in two equally sized subsamples: 2000.Q1-2009.Q4 and 2010.Q1-2019.Q4, of 40 observations each. As shown in Section 4, with such small sample sizes the DM framework with standard asymptotics suffers from large size distortions but fixed-smoothing asymptotics can still provide reliable inference.

Summary statistics for the test statistics  $d_t$  for unemployment, GDP and HCPI forecasts are in Tables 3-5. We report the full sample mean, the standard deviation, and autocorrelations up to the fourth lag and the eighth lag for the RPS and QPS loss functions.

Results for the equal predictive ability statistic are in the top panel of each table.

Table 3: Summary Statistics of  $d_t$  for unemployment rate

<b>Panel A: Equal Predictive Ability Statistic</b>							
<b>Loss</b>		<b>1 year ahead</b>			<b>2 years ahead</b>		
		<b>UG</b>	<b>GRW</b>	<b>Naïve</b>	<b>UG</b>	<b>GRW</b>	<b>Naïve</b>
RPS	<b>Mean</b>	1.24	0.24	0.25	0.70	0.21	0.21
	<b>STD</b>	1.15	0.37	0.40	1.57	0.63	0.48
	<b>AC1</b>	0.82	0.54	0.31	0.86	0.73	0.35
	<b>AC2</b>	0.60	0.20	-0.06	0.64	0.37	0.02
	<b>AC3</b>	0.36	0.00	-0.06	0.42	0.20	-0.04
	<b>AC4</b>	0.18	0.01	-0.06	0.24	0.16	-0.02
	<b>AC8</b>	-0.28	0.02	0.01	-0.24	-0.02	-0.07
QPS	<b>Mean</b>	0.16	0.05	0.10	0.06	0.00	0.05
	<b>STD</b>	0.24	0.22	0.18	0.23	0.16	0.14
	<b>AC1</b>	0.54	0.37	0.32	0.75	0.56	0.14
	<b>AC2</b>	0.33	0.31	-0.07	0.41	0.16	-0.08
	<b>AC3</b>	0.07	0.06	0.06	0.19	-0.06	-0.02
	<b>AC4</b>	0.02	0.07	0.06	0.08	-0.12	0.17
	<b>AC8</b>	-0.02	-0.06	0.17	-0.09	-0.32	0.06

<b>Panel B: Forecast Encompassing Statistic</b>							
<b>Loss</b>		<b>1 year ahead</b>			<b>2 years ahead</b>		
		<b>UG</b>	<b>GRW</b>	<b>Naïve</b>	<b>UG</b>	<b>GRW</b>	<b>Naïve</b>
RPS	<b>Mean</b>	0.05	0.01	-0.07	0.28	0.01	-0.07
	<b>STD</b>	0.48	0.16	0.16	0.87	0.28	0.19
	<b>AC1</b>	0.71	0.47	0.29	0.86	0.68	0.40
	<b>AC2</b>	0.44	0.01	-0.14	0.63	0.27	0.08
	<b>AC3</b>	0.28	0.07	-0.05	0.44	0.14	-0.05
	<b>AC4</b>	0.20	0.21	-0.10	0.27	0.17	-0.05
	<b>AC8</b>	-0.12	0.36	0.25	-0.26	0.27	-0.01
QPS	<b>Mean</b>	0.02	0.02	-0.03	0.05	0.03	-0.01
	<b>STD</b>	0.11	0.10	0.07	0.12	0.07	0.05
	<b>AC1</b>	0.47	0.35	0.32	0.76	0.55	0.13
	<b>AC2</b>	0.26	0.24	-0.11	0.40	0.07	-0.10
	<b>AC3</b>	0.02	0.02	0.10	0.18	-0.12	-0.08
	<b>AC4</b>	0.03	0.01	0.03	0.08	-0.11	0.08
	<b>AC8</b>	0.05	-0.08	0.25	-0.14	-0.27	0.10

Note: sample mean, standard deviation (STD), and autocorrelation coefficients up to order 4, and order 8 (AC1, AC2, AC3, AC4, AC8) of  $d_t$  for the unemployment rate. The top panel refers to  $d_t$  as defined for the equal predictive ability test in (3), and the bottom panel refers to  $d_t$  as defined for the forecast encompassing test in (6).

Table 4: Summary Statistics of  $d_t$  for GDP growth

<b>Panel A: Equal Predictive Ability Statistic</b>							
Loss		1 year ahead			2 years ahead		
		UG	GRW	Naïve	UG	GRW	Naïve
RPS	<b>Mean</b>	0.53	0.75	0.28	-0.28	0.59	0.03
	<b>STD</b>	1.21	1.70	0.60	0.86	1.95	0.25
	<b>AC1</b>	0.66	0.72	0.50	0.82	0.75	0.15
	<b>AC2</b>	0.32	0.45	0.18	0.60	0.34	-0.22
	<b>AC3</b>	0.09	0.28	0.18	0.44	0.02	-0.22
	<b>AC4</b>	-0.01	0.13	0.14	0.37	-0.10	-0.15
	<b>AC8</b>	-0.06	0.04	0.11	0.31	-0.05	0.11
QPS	<b>Mean</b>	0.04	0.09	0.04	-0.07	0.02	0.00
	<b>STD</b>	0.26	0.31	0.17	0.20	0.26	0.05
	<b>AC1</b>	0.48	0.47	0.17	0.75	0.72	0.17
	<b>AC2</b>	0.00	0.08	-0.04	0.45	0.37	-0.29
	<b>AC3</b>	-0.10	0.00	0.05	0.30	0.20	-0.22
	<b>AC4</b>	-0.11	0.01	-0.05	0.25	0.07	-0.01
	<b>AC8</b>	0.07	-0.06	0.03	0.19	-0.12	0.08
<b>Panel B: Forecast Encompassing Statistic</b>							
Loss		1 year ahead			2 years ahead		
		UG	GRW	Naïve	UG	GRW	Naïve
RPS	<b>Mean</b>	0.03	0.01	-0.08	0.30	0.17	-0.01
	<b>STD</b>	0.34	0.41	0.21	0.51	0.49	0.13
	<b>AC1</b>	0.51	0.71	0.45	0.84	0.79	0.12
	<b>AC2</b>	0.25	0.41	-0.01	0.60	0.49	-0.24
	<b>AC3</b>	0.04	0.15	-0.03	0.42	0.25	-0.23
	<b>AC4</b>	0.01	0.03	0.04	0.35	0.12	-0.17
	<b>AC8</b>	0.08	0.17	0.19	0.35	0.14	0.12
QPS	<b>Mean</b>	0.05	0.05	0.00	0.08	0.08	0.00
	<b>STD</b>	0.12	0.13	0.07	0.11	0.11	0.03
	<b>AC1</b>	0.49	0.48	0.11	0.77	0.76	0.13
	<b>AC2</b>	0.02	0.05	-0.13	0.49	0.49	-0.24
	<b>AC3</b>	-0.05	-0.06	0.03	0.34	0.34	-0.16
	<b>AC4</b>	-0.07	-0.10	-0.11	0.28	0.23	-0.04
	<b>AC8</b>	0.11	0.13	0.07	0.22	0.16	0.07

Note: sample mean, standard deviation (STD), and autocorrelation coefficients up to order 4, and order 8 (AC1, AC2, AC3, AC4, AC8) of  $d_t$  for GDP growth. The top panel refers to  $d_t$  as defined for the equal predictive ability test in (3), and the bottom panel refers to  $d_t$  as defined for the forecast encompassing test in (6).

Table 5: Summary Statistics of  $d_t$  for HCPI

<b>Panel A: Equal Predictive Ability Statistic</b>							
Loss		1 year ahead			2 years ahead		
		UG	GRW	Naïve	UG	GRW	Naïve
RPS	<b>Mean</b>	0.15	0.18	0.05	0.07	0.54	0.01
	<b>STD</b>	0.75	0.90	0.23	0.57	1.17	0.11
	<b>AC1</b>	0.71	0.54	0.16	0.75	0.65	-0.02
	<b>AC2</b>	0.60	0.33	-0.30	0.55	0.38	-0.02
	<b>AC3</b>	0.53	0.22	-0.03	0.38	0.25	0.03
	<b>AC4</b>	0.35	0.04	0.13	0.20	0.11	0.20
	<b>AC8</b>	0.03	-0.30	0.10	-0.01	-0.35	-0.09
QPS	<b>Mean</b>	-0.04	0.01	0.00	-0.02	0.08	0.00
	<b>STD</b>	0.21	0.25	0.09	0.18	0.31	0.06
	<b>AC1</b>	0.18	0.37	-0.14	0.35	0.60	-0.25
	<b>AC2</b>	0.17	0.19	-0.20	0.16	0.41	-0.12
	<b>AC3</b>	0.14	0.19	0.03	0.03	0.21	0.14
	<b>AC4</b>	-0.03	-0.03	-0.04	-0.24	-0.08	0.00
	<b>AC8</b>	-0.15	-0.20	0.01	-0.06	-0.19	-0.12

<b>Panel B: Forecast Encompassing Statistic</b>							
Loss		1 year ahead			2 years ahead		
		UG	GRW	Naïve	UG	GRW	Naïve
RPS	<b>Mean</b>	0.08	0.17	-0.02	0.07	0.03	0.00
	<b>STD</b>	0.33	0.49	0.11	0.28	0.42	0.05
	<b>AC1</b>	0.62	0.71	0.12	0.73	0.60	-0.02
	<b>AC2</b>	0.50	0.57	-0.31	0.52	0.43	-0.02
	<b>AC3</b>	0.44	0.32	-0.02	0.34	0.30	0.05
	<b>AC4</b>	0.26	0.11	0.13	0.15	0.08	0.19
	<b>AC8</b>	0.03	-0.21	0.08	-0.08	-0.26	-0.11
QPS	<b>Mean</b>	0.07	0.07	0.00	0.04	0.04	0.00
	<b>STD</b>	0.10	0.12	0.04	0.09	0.10	0.03
	<b>AC1</b>	0.19	0.34	-0.17	0.35	0.34	-0.26
	<b>AC2</b>	0.18	0.26	-0.17	0.14	0.20	-0.12
	<b>AC3</b>	0.14	0.14	0.02	0.02	0.04	0.16
	<b>AC4</b>	-0.03	0.01	-0.07	-0.26	-0.25	-0.01
	<b>AC8</b>	-0.15	-0.26	0.00	-0.07	-0.11	-0.12

Note: sample mean, standard deviation (STD), and autocorrelation coefficients up to order 4, and order 8 (AC1, AC2, AC3, AC4, AC8) of  $d_t$  for HCPI. The top panel refers to  $d_t$  as defined for the equal predictive ability test in (3), and the bottom panel refers to  $d_t$  as defined for the forecast encompassing test in (6).



A positive entry for the sample mean indicates that the forecasters are more accurate than the benchmark. We can see that this is usually the case, although occasionally the opposite also took place (most notably, for GDP at two years against the unconditional Gaussian benchmark; two instances have also been recorded for the QPS loss function for the HCPI, but for values very close to 0). The autocorrelation profile seems usually more relevant when the unconditional Gaussian benchmark is used. This may be worth noticing because in presence of large autocorrelation (relative to the sample size), the test may have low power, see Coroneo and Iacone (2021). Fortunately, in this case we can see that the autocorrelation coefficients at lag eight are usually quite small, indicating a quick decay of the dependence.

Results for the forecast encompassing statistic are in the bottom panel of each table. A negative entry for the sample mean is associated to an unrestricted negative estimate of the weight in the forecasts combination (which is not feasible given that the weight needs to be null or positive). We interpret this result as evidence that the SPF forecast cannot be improved through a linear combination with the benchmark. This happens in almost all cases for the naive benchmark. In Table 6, we report the estimated forecast combination weights for the full sample (from 2000.Q1 to 2019.Q4) and the two subsamples (2000.Q1 to 2009.Q4 and 2010.Q1 to 2019.Q4). The naive benchmark is usually associated with the smallest weight, the most notable exception being for real GDP growth forecast at two years ahead using the QPS loss function. Comparing across the two subsamples, we can see that the estimated forecast combination weights are similar across the two subsamples for the unemployment rate, but they are smaller in the second subsample for real GDP growth (and to a lesser extent for inflation), suggesting that the benchmarks contained less additional information about real GDP growth in the second part of the sample.

Density forecast evaluation test results are reported in Table 7, for the full sample, and in Tables 8 and 9, for the two subsamples. In Panel A, we report the equal predictive

Table 6: Estimated forecast combination weights

Full sample Q1.2000 - Q4.2019,  $T = 80$ .

Variable	Loss	1 year ahead			2 years ahead		
		UG	GRW	Naïve	UG	GRW	Naïve
UN	RPS	0.03	0.03	0.00	0.22	0.04	0.00
	QPS	0.11	0.24	0.00	0.33	0.49	0.00
GDP	RPS	0.05	0.01	0.00	0.92	0.18	0.00
	QPS	0.35	0.25	0.00	0.83	0.43	0.50
HCPI	RPS	0.26	0.32	0.00	0.33	0.05	0.00
	QPS	0.71	0.46	0.35	0.63	0.24	0.29

Subsample Q1.2000 - Q4.2009,  $T = 40$ .

Variable	Loss	1 year ahead			2 years ahead		
		UG	GRW	Naïve	UG	GRW	Naïve
UN	RPS	0.06	0.09	0.00	0.33	0.23	0.00
	QPS	0.11	0.20	0.00	0.34	0.45	0.00
GDP	RPS	0.12	0.06	0.00	1.00	0.17	0.00
	QPS	0.47	0.35	0.08	1.00	0.61	0.29
HCPI	RPS	0.58	0.24	0.00	0.70	0.07	0.76
	QPS	0.76	0.34	0.79	0.65	0.20	0.93

Subsample Q1.2010 - Q4.2019,  $T = 40$ .

Variable	Loss	1 year ahead			2 years ahead		
		UG	GRW	Naïve	UG	GRW	Naïve
UN	RPS	0.01	0.00	0.00	0.07	0.00	0.00
	QPS	0.12	0.33	0.00	0.30	0.58	0.00
GDP	RPS	0.00	0.00	0.00	0.45	0.20	0.00
	QPS	0.20	0.15	0.00	0.48	0.24	0.95
HCPI	RPS	0.00	0.50	0.00	0.00	0.01	0.00
	QPS	0.64	0.81	0.00	0.58	0.34	0.00

Note: LS estimate of  $\lambda$  in (5). The top panel refers to the full sample (Q1.2000 - Q4.2019,  $T = 80$ ), the middle panel to the first half-sample (Q1.2000 - Q4.2009,  $T = 40$ ) and the bottom panel to the second-half sample (Q1.2010 - Q4.2019,  $T = 40$ ).

accuracy test, and in Panel B we report the forecast encompassing test. A negative value of the equal predictive accuracy test indicates that the benchmark is performing better than the ECB SPF forecast, while a negative value for the forecast encompassing test indicates that the unrestricted weight on the benchmark is negative, as it does not have any additional information with respect to the ECB SPF forecast. Rejections

from standard asymptotics critical values are indicated shading the appropriate cell;  $\blacksquare$  and  $\square$  indicate, respectively, one-sided significance at the 5% and 10% level. Rejections using fixed-smoothing asymptotics critical values are reported using \*\* and \* to indicate, respectively, one-sided significance at the 5% and 10% level.

For each variable and each test (equal predictive accuracy and forecast encompassing), we report results with both the Ranked Probability Score (RPS) and the Quadratic Probability Score (QPS) loss functions. The long run variance is estimated using WCE with the Bartlett kernel and bandwidth  $M = \lfloor T^{1/2} \rfloor$  and WPE with Daniell kernel and bandwidth  $m = \lfloor T^{1/3} \rfloor$ .

For the unemployment rate, results for the equal predictive ability test indicate that one-year ahead ECB SPF forecasts outperforms the unconditional Gaussian and the naive benchmarks in all three samples; while they outperform the Gaussian random walk only according to the RPS loss function. For the two-years ahead forecasts, the surveys again beat the unconditional Gaussian benchmark, and the ECB SPF forecasts superiority is more marked in the second half of the sample. As for the forecast encompassing test, the ECB SPF encompasses all the benchmarks at all horizons, with the RPS loss function; when using the QPS loss function we find that the two Gaussian benchmarks are not encompassed by the ECB SPF. Overall, these results indicate that for unemployment the ECB SPF provides more accurate one year-ahead predictions than the benchmarks, however for two-year ahead predictions the ECB SPF forecasts can be improved by combining them with the Gaussian benchmarks.

Results for real GDP growth are less favourable for the ECB SPF, especially when using fixed-smoothing asymptotics. At one-year ahead, there is not conclusive evidence that it delivers more accurate predictions than the benchmarks in the first subsample. On the other hand, in the second subsample, at one year ahead the null of equal forecast accuracy using the RPS loss is rejected for all benchmarks. At two years ahead, we cannot reject the null of equal forecast accuracy of the ECB SPF forecasts and the

Table 7: Forecast evaluation tests. Full sample Q1.2000 - Q4.2019,  $T = 80$ .

<b>Panel A: Equal Predictive Ability Test</b>								
Variable	Loss	LRV	1 year ahead			2 years ahead		
			UG	GRW	Naïve	UG	GRW	Naïve
UN	RPS	WCE	4.95**	3.95**	5.42**	1.98**	1.68*	3.74**
		WPE	4.17**	3.63**	6.58**	1.71*	1.53*	4.99**
	QPS	WCE	3.89**	1.47*	4.22**	1.25	0.05	2.97**
		WPE	3.68**	1.37	4.13**	1.10	0.05	3.46**
GDP	RPS	WCE	2.39**	2.13**	2.70**	-1.38	1.68*	1.47*
		WPE	2.22**	1.79*	2.44**	-1.12	1.52*	1.57*
	QPS	WCE	1.07	1.99**	1.91*	-1.51	0.48	0.02
		WPE	0.88	1.59*	1.72*	-1.27	0.43	0.02
HCPI	RPS	WCE	0.87	1.13	2.19**	0.56	2.40**	1.14
		WPE	0.74	1.09	1.74*	0.50	2.47**	0.99
	QPS	WCE	-1.41	0.33	0.40	-0.70	1.51*	0.31
		WPE	-1.40	0.38	0.37	-0.71	1.51*	0.27

<b>Panel B: Forecast Encompassing Test</b>								
Variable	Loss	LRV	1 year ahead			2 years ahead		
			UG	GRW	Naïve	UG	GRW	Naïve
UN	RPS	WCE	0.46	0.24	-4.18	1.39	0.18	-2.83
		WPE	0.36	0.21	-4.24	1.11	0.16	-3.53
	QPS	WCE	1.24	1.58*	-2.68	2.30**	3.13**	-1.74
		WPE	1.16	1.49*	-2.40	1.96**	3.06**	-1.89
GDP	RPS	WCE	0.53	0.11	-2.63	2.60**	1.59*	-0.52
		WPE	0.47	0.10	-2.46	2.09**	1.40*	-0.60
	QPS	WCE	2.46**	2.50**	-0.04	3.48**	3.27**	1.47*
		WPE	2.03**	2.01**	-0.03	2.86**	2.76**	1.54*
HCPI	RPS	WCE	1.12	1.71*	-1.47	1.11	0.35	-0.57
		WPE	0.97	1.43*	-1.18	1.00	0.33	-0.50
	QPS	WCE	4.65**	3.97**	0.95	3.59**	2.71**	0.43
		WPE	4.75**	3.44**	0.92	3.91**	2.59**	0.37

Note: Equal predictive ability test statistics and the forecast encompassing test statistics for one-year and two-year ahead ECB SPF density forecasts against the unconditional Gaussian, the Gaussian random walk and the naive benchmark forecasts on the full sample Q1.2000 - Q4.2019 ( $T = 80$ ). A negative equal predictive ability test statistic sign implies that benchmark performs better than the ECB SPF, and a negative value for the forecast encompassing test indicates that the estimated unrestricted weight on the benchmark is negative. Long run variances are estimated using WCE with Bartlett kernel and bandwidth  $M = \lfloor T^{1/2} \rfloor$  and WPE with Daniell kernel and bandwidth  $m = \lfloor T^{1/3} \rfloor$ . ■ and ■ indicate, respectively, one-sided significance at the 5% and 10% level using standard asymptotics. Rejections using fixed-smoothing asymptotics are reported using \*\* and \* to indicate, respectively, one-sided significance at the 5% and 10% level.

Table 8: Forecast evaluation tests. Subsample Q1.2000 - Q4.2009,  $T = 40$ .

<b>Panel A: Equal Predictive Ability Test</b>								
Variable	Loss	LRV	1 year ahead			2 years ahead		
			UG	GRW	Naïve	UG	GRW	Naïve
UN	RPS	WCE	3.52**	3.22**	2.80**	0.81	1.25	2.27**
		WPE	2.98**	2.90**	2.37**	0.71	1.25	1.87*
	QPS	WCE	2.69**	1.34	2.15**	0.89	0.20	2.31**
		WPE	2.25**	1.08	1.93*	0.78	0.16	2.15**
GDP	RPS	WCE	1.34	1.61*	2.20**	-2.21	1.24	1.04
		WPE	1.10	1.46*	2.11**	-1.98	1.12	1.31
	QPS	WCE	0.17	1.00	1.36	-4.24	-1.05	0.63
		WPE	0.14	0.81	1.35	-4.07	-0.92	0.62
HCPI	RPS	WCE	-0.46	1.62*	0.90	-0.83	1.85*	-0.27
		WPE	-0.38	1.47*	0.83	-0.71	1.53*	-0.26
	QPS	WCE	-1.47	1.40	-0.54	-0.71	1.56*	-0.49
		WPE	-1.30	1.11	-0.49	-0.58	1.51*	-0.48

<b>Panel B: Forecast Encompassing Test</b>								
Variable	Loss	LRV	1 year ahead			2 years ahead		
			UG	GRW	Naïve	UG	GRW	Naïve
UN	RPS	WCE	0.50	1.32	-1.74	1.46	1.41	-1.70
		WPE	0.43	1.06	-1.93	1.24	1.63*	-1.41
	QPS	WCE	0.73	0.99	-0.89	1.82*	2.08**	-1.23
		WPE	0.61	0.81	-0.86	1.54*	1.76*	-1.15
GDP	RPS	WCE	0.91	0.39	-1.99	2.76**	1.09	-0.13
		WPE	0.75	0.35	-2.05	2.48**	1.00	-0.17
	QPS	WCE	2.73**	2.61**	0.35	5.50**	4.86**	0.75
		WPE	2.45**	2.30**	0.36	5.43**	4.90**	0.77
HCPI	RPS	WCE	4.08**	1.04	-0.30	2.76**	0.52	0.69
		WPE	3.43**	0.96	-0.29	2.48**	0.49	0.67
	QPS	WCE	4.13**	2.98**	1.58*	3.33**	2.78**	1.00
		WPE	3.62**	2.97**	1.48*	2.65**	2.17**	1.01

Note: Equal predictive ability test statistics and the forecast encompassing test statistics for one-year and two-year ahead ECB SPF density forecasts against the unconditional Gaussian, the Gaussian random walk and the naive benchmark forecasts on the subsample Q1.2000 - Q4.2009 ( $T = 40$ ). A negative equal predictive ability sign implies that the benchmark performs better than the ECB SPF, and a negative value for the forecast encompassing test indicates that the unrestricted weight on the benchmark is negative. Long run variances are estimated using WCE with Bartlett kernel and bandwidth  $M = \lfloor T^{1/2} \rfloor$  and WPE with Daniell kernel and bandwidth  $m = \lfloor T^{1/3} \rfloor$ . ■ and ■ indicate, respectively, one-sided significance at the 5% and 10% level using standard asymptotics. Rejections using fixed-smoothing asymptotics are reported using \*\* and \* to indicate, respectively, one-sided significance at the 5% and 10% level.

Table 9: Forecast evaluation tests. Subsample Q1.2010 - Q4.2019,  $T = 40$ .

<b>Panel A: Equal Predictive Ability Test</b>								
Variable	Loss	LRV	1 year ahead			2 years ahead		
			UG	GRW	Naïve	UG	GRW	Naïve
UN	RPS	WCE	3.80**	2.23**	5.10**	2.71**	1.35	2.34**
		WPE	3.17**	2.00**	4.82**	2.39**	1.26	2.10**
	QPS	WCE	3.30**	0.73	5.29**	0.92	-0.32	2.40**
		WPE	2.77**	0.66	5.46**	0.79	-0.30	2.18**
GDP	RPS	WCE	3.19**	2.82**	2.67**	0.14	1.80*	0.86
		WPE	2.69**	3.35**	2.50**	0.13	1.97**	0.73
	QPS	WCE	1.58*	2.01**	1.91*	0.07	1.16	-0.61
		WPE	1.49*	2.10**	1.53*	0.06	1.20	-0.58
HCPI	RPS	WCE	1.23	-0.01	2.32**	1.34	1.62*	1.98**
		WPE	1.10	-0.01	1.90*	1.15	1.25	1.97**
	QPS	WCE	-0.49	-0.92	1.80*	-0.24	0.47	1.40
		WPE	-0.41	-0.74	2.03**	-0.21	0.38	1.54*

<b>Panel B: Forecast Encompassing Test</b>								
Variable	Loss	LRV	1 year ahead			2 years ahead		
			UG	GRW	Naïve	UG	GRW	Naïve
UN	RPS	WCE	0.15	-0.51	-4.80	0.46	-0.50	-1.84
		WPE	0.12	-0.45	-4.79	0.41	-0.46	-1.67
	QPS	WCE	1.63*	1.67*	-4.18	1.64*	2.58**	-1.41
		WPE	1.44*	1.53*	-4.32	1.45*	2.34**	-1.31
GDP	RPS	WCE	-0.50	-0.71	-2.25	1.35	1.54*	-0.53
		WPE	-0.42	-1.02	-2.12	1.28	1.25	-0.44
	QPS	WCE	1.04	1.15	-0.91	1.39	1.35	1.33
		WPE	0.97	1.14	-0.73	1.30	1.24	1.26
HCPI	RPS	WCE	-0.08	1.72*	-1.76	-0.33	0.03	-1.63
		WPE	-0.07	1.39	-1.41	-0.28	0.02	-1.65
	QPS	WCE	2.41**	2.59**	-0.66	2.00**	1.22	-0.82
		WPE	1.96**	2.10**	-0.72	1.78*	0.99	-0.91

Note: Equal predictive ability test statistic and forecast encompassing test statistics values for one-year and two-year ahead ECB SPF density forecasts against the unconditional Gaussian, the Gaussian random walk and the naive benchmark forecasts on the subsample Q1.2010 - Q4.2019 ( $T = 40$ ). A negative value for the equal predictive ability test indicates that benchmarks perform better than the ECB SPF, and a negative value for the forecast encompassing test indicates that the unrestricted weight on the benchmark is negative. Long run variances are estimated using WCE with Bartlett kernel and bandwidth  $M = \lfloor T^{1/2} \rfloor$  and WPE with Daniell kernel and bandwidth  $m = \lfloor T^{1/3} \rfloor$ . ■ and ■ indicate, respectively, one-sided significance at the 5% and 10% level using standard asymptotics. Rejections using fixed-smoothing asymptotics are reported using \*\* and \* to indicate, respectively, one-sided significance at the 5% and 10% level.

benchmarks. This is confirmed by the forecast encompassing test, that indicates that at two-years ahead the Gaussian benchmarks are not encompassed by the ECB SPF in the full sample and the first subsample. Overall, these results indicate that for real GDP growth, ECB SPF forecasters can outperform simple benchmarks at least at one-year horizon, especially in the second subsample.

In the case of inflation, there is no statistical evidence that ECB SPF density forecasts outperform the benchmarks from the equal predictive ability test, as the null hypothesis of equal forecast accuracy is almost never rejected. However, in the second subsample the ECB SPF outperforms the naive benchmark. The forecast encompassing test results indicate that the ECB SPF density forecast for inflation do not encompass the Gaussian benchmarks, suggesting that more accurate density forecasts for inflation can be obtained by combining these benchmarks with the ECB SPF.

Overall, using the QPS yields qualitatively the same outcome as using the RPS, although significant results for the equal predictive ability (forecast rationality) tests are less (more) frequent when the QPS is used, see for example the case of the real GDP growth. For the forecast encompassing test, we find that the null of no encompassing is rejected more often with the QPS, indicating that the ECB SPF place more probability in the neighbourhood of the effective outcome, often near-missing the true realization.

As for the benchmarks, the ECB SPF easily outperforms and encompasses the naive benchmark, indicating that the professional forecasters update their information set when making their predictions and that previous round forecasts are completely uninformative. On the other hand, the unconditional Gaussian benchmark seems the most difficult to outperform and encompass, especially for two-years ahead forecasts.

Comparing the application of standard asymptotics with fixed-smoothing asymptotics, we reject the null of equal predictive ability more frequently for the tests with standard asymptotics, especially when the sample is split in two parts, and for long-horizon forecasts. This is due to the fact that in the subsamples the tests are performed

only on 40 observations, exacerbating the size distortions induced by standard asymptotics, see Section 4. For example, in Table 9 both tests with standard asymptotics reject at 10% significance level the null of no encompassing for GDP growth, at the two-year ahead horizon. This could be interpreted as indication that the benchmarks are not encompassed by the SPF. However, these results are mostly not confirmed when using fixed-smoothing asymptotics, indicating that they are partially spurious and demonstrating the risks of using standard asymptotics in a small sample.

## 6 Conclusions

In this paper, we apply fixed- $b$  and fixed- $m$  asymptotics to tests of equal predictive accuracy and encompassing for survey density forecasts. In an original Monte Carlo design, we verify that fixed-smoothing asymptotics delivers correctly sized tests in this framework, even when only a small number of out of sample observations is available.

We apply the density forecast evaluation tests with fixed-smoothing asymptotics to evaluate the predictive ability of density forecasts from the European Central Bank's Survey of Professional Forecasters (ECB SPF) over the period 2001.Q1-2019.Q4, taking as benchmarks simple forecasts generated from an unconditional Gaussian distributions, a Gaussian random walk and the previous survey round.

Our results indicate that ECB SPF density forecasts for unemployment and real GDP growth outperformed and sometimes encompassed the benchmarks, especially at one-year ahead and in the second subsample. On the contrary, survey forecasts for inflation do not easily outperform nor encompass the benchmarks. For all the variables, however, we find evidence of an improvement in predictive ability since 2010, supporting the anecdotal evidence of a change in the forecasting practice after the financial crisis.



## A Density forecast encompassing test

In this Appendix, we show the null hypothesis of density forecast encompassing can be tested using the DM framework by defining  $d_t$  as in (6).

If we denote the forecast errors associated with  $\mathbf{f}_{t,c}(\lambda)$  in (5) as  $\mathbf{e}_{t,c}(\lambda) = \mathbf{y}_t - \mathbf{f}_{t,c}(\lambda)$ , then, the optimal weight in the minimum QPS sense has

$$\hat{\lambda} = \arg \min \sum_{t=1}^T (\mathbf{e}_{t,c}(\lambda))' \mathbf{e}_{t,c}(\lambda). \quad (12)$$

The derivative is

$$\frac{\partial}{\partial \lambda} \sum_{t=1}^T (\mathbf{e}_{t,c}(\lambda))' \mathbf{e}_{t,c}(\lambda) = \sum_{t=1}^T 2(\mathbf{y}_t - \mathbf{f}_{t,c}(\lambda))' \frac{\partial}{\partial \lambda} (-\mathbf{f}_{t,c}(\lambda)) = \sum_{t=1}^T 2(\mathbf{y}_t - \mathbf{f}_{t,c}(\lambda))' (\mathbf{f}_{t,1} - \mathbf{f}_{t,2})$$

and the first order condition therefore gives

$$\sum_{t=1}^T 2(\mathbf{y}_t - \mathbf{f}_{t,c}(\lambda))' (\mathbf{f}_{t,1} - \mathbf{f}_{t,2}) = 0.$$

which is met for  $\lambda = \hat{\lambda}$  (i.e.,  $\hat{\lambda}$  is defined in this way).

Let

$$d_t(\lambda) = -(\mathbf{y}_t - \mathbf{f}_{t,c}(\lambda))' (\mathbf{f}_{t,1} - \mathbf{f}_{t,2})$$

If  $\mathbf{y}_t$ ,  $\mathbf{f}_{1,t}$  and  $\mathbf{f}_{2,t}$  are jointly mixing with a sufficient rate, then so is  $d_t(\lambda)$ .

Denoting  $\sigma_T^2(\lambda) = \text{Var}(\sqrt{T} \frac{1}{T} \sum_{t=1}^T d_t(\lambda))$  as the long run variance, assuming that  $d_t(\lambda)$  is mixing with sufficient rate and  $\sigma_T(\lambda) > 0$  then we have a CLT for standardised sum of  $d_t(\lambda)$ . This suggests a LM type test for forecast encompassing. Mimicking the first order condition, denote  $\lambda_0$  as the value of  $\lambda$  that gives  $E(d_t(\lambda))|_{\lambda=\lambda_0} = 0$ , then

$$\sqrt{T} \frac{1/T \sum_{t=1}^T d_t(\lambda_0)}{\sigma_T(\lambda_0)} \rightarrow_d N(0, 1)$$

and, under  $H_0 : \{\lambda_0 = 0\}$  then

$$\sqrt{T} \frac{1/T \sum_{t=1}^T d_t}{\sigma_T} \rightarrow_d N(0, 1)$$

where we used  $d_t$  and  $\sigma_T$  in place of  $d_t(0)$  and  $\sigma_T(0)$  to shorten the notation.

Rewriting

$$d_t(\lambda) = -(\mathbf{y}_t - \mathbf{f}_{t,1} - \lambda(\mathbf{f}_{t,1} - \mathbf{f}_{t,2}))'(\mathbf{f}_{t,1} - \mathbf{f}_{t,2})$$

then

$$d_t = \mathbf{e}'_{t,1}(\mathbf{e}_{t,1} - \mathbf{e}_{t,2})$$

These facts suggest for  $H_0 : \{\lambda_0 = 0\}$  the test statistic

$$\sqrt{T} \frac{1/T \sum_{t=1}^T d_t}{\hat{\sigma}} = \sqrt{T} \frac{1/T \sum_{t=1}^T \mathbf{e}'_{t,1}(\mathbf{e}_{t,1} - \mathbf{e}_{t,2})}{\hat{\sigma}}$$

for an appropriate estimate of the long run variance  $\hat{\sigma}$ .

To complete the specification of the test and to check the power, we rewrite

$$\begin{aligned} \frac{\sqrt{T}}{T} \sum_{t=1}^T d_t &= \frac{\sqrt{T}}{T} \sum_{t=1}^T (d_t - d_t(\lambda_0) + d_t(\lambda_0)) \\ &= \frac{\sqrt{T}}{T} \sum_{t=1}^T (d_t - d_t(\lambda_0)) + \frac{\sqrt{T}}{T} \sum_{t=1}^T (d_t(\lambda_0)) \\ &= \frac{\sqrt{T}}{T} \sum_{t=1}^T (d_t - d_t(\lambda_0)) + O_p(1) \end{aligned}$$

and notice that

$$\begin{aligned} d_t - d_t(\lambda_0) &= -(\mathbf{y}_t - \mathbf{f}_{t,c}(0))'(\mathbf{f}_{t,1} - \mathbf{f}_{t,2}) + (\mathbf{y}_t - \mathbf{f}_{t,c}(\lambda_0))'(\mathbf{f}_{t,1} - \mathbf{f}_{t,2}) \\ &= (\mathbf{f}_{t,c}(0) - \mathbf{f}_{t,c}(\lambda_0))'(\mathbf{f}_{t,1} - \mathbf{f}_{t,2}) \\ &= (\mathbf{f}_{t,1} - (1 - \lambda_0)\mathbf{f}_{t,1} - \lambda_0\mathbf{f}_{t,2})'(\mathbf{f}_{t,1} - \mathbf{f}_{t,2}) \\ &= \lambda_0(\mathbf{f}_{t,1} - \mathbf{f}_{t,2})'(\mathbf{f}_{t,1} - \mathbf{f}_{t,2}) = \lambda_0(\mathbf{e}_{t,1} - \mathbf{e}_{t,2})'(\mathbf{e}_{t,1} - \mathbf{e}_{t,2}) \end{aligned}$$

Thus, for the alternative  $H_A : \{\lambda_0 > 0\}$  the null hypothesis is rejected if the test statistic takes a value larger than the critical value.

Notice that the value that solves  $E(d_t(\lambda_0)) = 0$  is

$$\lambda_0 = \frac{E(\mathbf{e}'_{t,1}(\mathbf{e}_{t,1} - \mathbf{e}_{t,2}))}{E(\mathbf{e}_{t,1} - \mathbf{e}_{t,2})'(\mathbf{e}_{t,1} - \mathbf{e}_{t,2})}$$

so if, for example,  $\mathbf{e}_{t,1}$  and  $\mathbf{e}_{t,2}$  are vectors of independent, identically distributed sequences, independent from each other, then  $\lambda_0 = 1/2$ . On the other hand, if  $E(\mathbf{e}'_{t,1}(\mathbf{e}_{t,1} - \mathbf{e}_{t,2})) = 0$  then  $\lambda_0 = 0$ .

The Ranked Probability Score (RPS) may be treated in the same way, using the cumulative distribution functions of each density forecast  $\mathbf{F}_{t,i}$  and of the individual realisation  $\mathbf{Y}_t$ .

## References

- Andrews, Donald W. K. (1984) ‘Non-strong mixing autoregressive processes.’ *Journal of Applied Probability* 21(4), 930–934
- Boero, Gianna, Jeremy Smith, and Kenneth F Wallis (2008) ‘Uncertainty and disagreement in economic prediction: the Bank of England Survey of External Forecasters.’ *The Economic Journal* 118(530), 1107–1127
- Bowles, Carlos, Roberta Friz, Veronique Genre, Geoff Kenny, Aidan Meyler, and Tuomas Rautanen (2007) ‘The ECB Survey of Professional Forecasters (SPF)-a review after eight years’ experience.’ *ECB Occasional Paper*
- Brier, Glenn W (1950) ‘Verification of forecasts expressed in terms of probability.’ *Monthly Weather Review* 78(1), 1–3
- Choi, Hwan-sik, and Nicholas M Kiefer (2010) ‘Improving robust model selection tests for dynamic models.’ *The Econometrics Journal* 13(2), 177–204
- Clark, Todd E (1999) ‘Finite-sample properties of tests for equal forecast accuracy.’ *Journal of Forecasting* 18(7), 489–504
- Clements, Michael P (2014) ‘Forecast uncertainty—ex ante and ex post: US inflation and output growth.’ *Journal of Business and Economic Statistics* 32(2), 206–216
- Clements, Michael P, and David I Harvey (2010) ‘Forecast encompassing tests and probability forecasts.’ *Journal of Applied Econometrics* 25(6), 1028–1062
- Coroneo, Laura, and Fabrizio Iacone (2020) ‘Comparing predictive accuracy in small samples using fixed-smoothing asymptotics.’ *Journal of Applied Econometrics* 35(4), 391–409
- (2021) ‘Testing for equal predictive accuracy with strong dependence.’ Technical Report, Department of Economics, University of York
- Coroneo, Laura, Fabrizio Iacone, Alessia Paccagnini, and Paulo Santos Monteiro (2022) ‘Testing the predictive accuracy of covid-19 forecasts.’ *International Journal of Forecasting*
- D’Amico, Stefania, Athanasios Orphanides et al. (2008) ‘Uncertainty and disagreement in economic forecasting.’ *Federal Reserve Board Finance and Economics Discussion Series*
- Dawid, A Philip (1984) ‘Present position and potential developments: Some personal views statistical theory the prequential approach.’ *Journal of the Royal Statistical Society: Series A (General)* 147(2), 278–290
- de Vincent-Humphreys, Rupert, Ivelina Dimitrova, Elisabeth Falck, and Lukas Henkel (2019) ‘Twenty years of the ECB survey of professional forecasters.’ *Economic Bulletin Articles*

- Diebold, Francis X, and Robert S Mariano (1995) ‘Comparing predictive accuracy.’ *Journal of Business and Economic Statistics* 13(3), 253–262
- Diebold, Francis X, Anthony S Tay, and Kenneth F Wallis (1999) ‘Evaluating density forecasts of inflation: the survey of professional forecasters.’ in R. Engle and H. White (eds.), *A Festschrift in Honour of Clive WJ Granger* pp. 76–90
- Epstein, Edward S (1969) ‘A scoring system for probability forecasts of ranked categories.’ *Journal of Applied Meteorology* 8(6), 985–987
- Fair, Ray C (1980) ‘Estimating the expected predictive accuracy of econometric models.’ *International Economic Review* 21(2), 355–378
- Garcia, Juan A (2003) ‘An introduction to the ECB’s Survey of Professional Forecasters.’ *ECB Occasional Paper*
- Giacomini, Raffaella, and Halbert White (2006) ‘Tests of conditional predictive ability.’ *Econometrica* 74(6), 1545–1578
- Giannone, Domenico, Jérôme Henry, Magdalena Lalik, and Michele Modugno (2012) ‘An area-wide real-time database for the Euro area.’ *Review of Economics and Statistics* 94(4), 1000–1013
- Gneiting, Tilmann, and Adrian E Raftery (2007) ‘Strictly proper scoring rules, prediction, and estimation.’ *Journal of the American Statistical Association* 102(477), 359–378
- Harvey, David I., Stephen J. Leybourne, and Emily J. Whitehouse (2017) ‘Forecast evaluation tests and negative long-run variance estimates in small samples.’ *International Journal of Forecasting* 33(4), 833 – 847
- Harvey, David I, Stephen J Leybourne, and Paul Newbold (1998) ‘Tests for forecast encompassing.’ *Journal of Business and Economic Statistics* 16(2), 254–259
- Hualde, Javier, and Fabrizio Iacone (2017) ‘Fixed bandwidth asymptotics for the studentized mean of fractionally integrated processes.’ *Economics Letters* 150, 39–43
- Kiefer, Nicholas M, and Timothy J Vogelsang (2002a) ‘Heteroskedasticity–autocorrelation robust standard errors using the bartlett kernel without truncation.’ *Econometrica* 70(5), 2093–2095
- (2002b) ‘Heteroskedasticity-autocorrelation robust testing using bandwidth equal to sample size.’ *Econometric Theory* 18(6), 1350–1366
- (2005) ‘A new asymptotic theory for heteroskedasticity-autocorrelation robust tests.’ *Econometric Theory* 21(6), 1130–1164
- Lazarus, Eben, Daniel J. Lewis, James H. Stock, and Mark W. Watson (2018) ‘HAR inference: Recommendations for practice.’ *Journal of Business and Economic Statistics* 36(4), 541–559

- Li, Jia, and Andrew J. Patton (2018) ‘Asymptotic inference about predictive accuracy using high frequency data.’ *Journal of Econometrics* 203(2), 223 – 240
- Manzan, Sebastiano (2021) ‘Are professional forecasters bayesian?’ *Journal of Economic Dynamics and Control* 123, 104045
- Mitchell, James, and Stephen G Hall (2005) ‘Evaluating, comparing and combining density forecasts using the klic with an application to the bank of england and niers fan charts of inflation.’ *Oxford Bulletin of Economics and Statistics* 67, 995–1033
- Neave, Henry R (1970) ‘An improved formula for the asymptotic variance of spectrum estimates.’ *The Annals of Mathematical Statistics* 41(1), 70–77
- Newey, Whitney K., and Kenneth D. West (1994) ‘Automatic lag selection in covariance matrix estimation.’ *The Review of Economic Studies* 61(4), 631–653
- Phillips, Peter C. B. (2005) ‘HAC estimation by automated regression.’ *Econometric Theory* 21(1), 116–142
- Phillips, Peter C. B., and Victor Solo (1992) ‘Asymptotics for linear processes.’ *Annals of Statistics* 20(2), 971–1001
- Sun, Yixiao (2013) ‘A heteroskedasticity and autocorrelation robust F test using an orthonormal series variance estimator.’ *The Econometrics Journal* 16(1), 1–26
- Sun, Yixiao (2014) ‘Let’s fix it: Fixed-b asymptotics versus small-b asymptotics in heteroskedasticity and autocorrelation robust inference.’ *Journal of Econometrics* 178, 659 – 677
- Tay, Anthony S, and Kenneth F Wallis (2000) ‘Density forecasting: a survey.’ *Journal of Forecasting* 19(4), 235–254
- Wooldridge, Jeffrey M., and Halbert White (1988) ‘Some invariance principles and central limit theorems for dependent heterogeneous processes.’ *Econometric Theory* 4(2), 210–230