

All Forecasters Are Not the Same: Systematic Patterns in Predictive Performance

Robert Rich
Federal Reserve Bank of Cleveland

Joseph Tracy
Daniels School of Business, Purdue University
American Enterprise Institute

October, 2024

Abstract: Are all forecasters the same? Expectations models incorporating information rigidities typically imply forecasters are interchangeable which predicts an absence of systematic patterns in individual forecast behavior. Motivated by this prediction, we examine the European Central Bank's Survey of Professional Forecasters and find, in contrast, that participants display systematic patterns in predictive performance both within and across target variables. Moreover, we document a new result from professional forecast surveys which is that inter- and intra-forecaster relative predictive performance are strongly linked to the degree of difficulty in the forecasting environment. This insight can inform the ongoing development of expectations models.

Keywords: forecasting profession; surveys; evaluating forecasts; point forecasts; density forecasts; heterogeneity

I. Introduction

Expectations are important for understanding the decision-making of households and firms, as well as for explaining movements in economic and financial variables. Early work on the formation of beliefs posited that agents form their expectations in a static or adaptive manner. However, these models eventually drew criticism because of their restrictions on agents' information sets and for allowing agents to make systematic forecast errors. In response, the full-information rational expectations (FIRE) model was developed which assumes that all agents know the true structure of the economy and have access to the same information set.

While the FIRE model remains the main paradigm for the formation of expectations, it implies that agents display identical forecast behavior and, therefore, cannot generate the type of dispersion in agents' expectations – that is, disagreement – observed in surveys or financial markets. Consequently, in recent years the FIRE model has been replaced with a weaker form of rational expectations in which agents use available information efficiently subject to certain constraints. A prominent feature of these models is the presence of informational rigidities (IR) either in the form of sticky information [Mankiw and Reis (2002); Mankiw, Reis, and Wolfers (2003)] or noisy information [Woodford (2003); Sims (2003); Mackowiak and Wiederholt (2009)].

While IR models can generate disagreement, a key, but largely overlooked, implication of almost all these models is that heterogeneity in individual forecast behavior should not display systematic patterns.¹ This is because variation in forecast behavior only arises from randomness either in the updating of individual information sets or in the configuration of shocks faced by individuals. Agents can display differences in their forecast behavior at a point in time, but their forecast behavior should be the same on average over time. Consequently, forecasters should be viewed as interchangeable with no distinguishing patterns in their average observed behavior.²

Motivated by this consideration, this study uses data from the European Central Bank's Survey of Professional Forecasters (ECB-SPF) to explore the implications of interchangeability across three aspects of forecast behavior. The first is whether the predictive performance metrics of participants display “distributional homogeneity” within a target variable – that is, do the mean and

¹ Coibion and Gorodnichenko (2012, 2015) test the predictions of the sticky information and the noisy information models for various aspects of forecast behavior at the aggregate level, but they do not consider this implication of IR models at the individual level.

² Clements (2022) makes this same observation to motivate his analysis of the US Survey of Professional Forecasters. The average observed behavior of forecasters refers to a period long enough to allow individuals to update their information sets on a comparable basis in the case of the sticky information model, or to be subject to a comparable set of shocks in the case of the noisy information model.

variation over time in a participant's accuracy align with those of others on average.³ The second is whether the relative accuracy of forecasters displays systematic patterns within a target variable. The third is whether individual forecasters display similar behavior across target variables. While these empirical features are of general interest for expectations models, we have noted their relevance for IR models.

We conduct the analysis within a panel data framework using the common correlated effects (CCE) estimator of Pesaran (2006). The CCE modeling strategy is attractive for several reasons. First, it allows for a broad identification of heterogeneity and correlation patterns in participants' predictive performance. Second, the inclusion in the individual regressions of a time-specific cross-sectional average of predictive performance controls for aggregate shocks that can generate dependence across participants which is a typical concern in panel data models. Last, the average predictive performance variable also provides a natural basis to identify tranquil/volatile forecast episodes that play a critical role in the analysis.

The results provide strong evidence that ECB-SPF participants are not interchangeable. Our tests reject the property of distributional homogeneity which indicates there are significant differences across forecasters in the mean and variance of their predictive performance metrics within a target variable. There are also systematic patterns in participants' relative predictive performance over time, both within and across target variables. Moreover, we document a new finding that these systematic patterns are strongly linked to the degree of difficulty in the forecasting environment. Within a target variable, some participants display higher relative accuracy in tranquil episodes, while other participants display higher relative accuracy in volatile episodes. Across target variables, we find that participants who display higher (lower) relative accuracy in tranquil/volatile environments for one target variable tend to display the same behavior for the other target variables. These results pose a challenge to IR models.

Our results are consistent with recent work by Hounyo and Lahiri (2023) who test for equal predictive ability among participants in the US Survey of Professional Forecasters (US-SPF) and report evidence of "persistent performance heterogeneity". Hounyo and Lahiri (2023) improve upon the bootstrap technique of D'Agostino, McQuinn and Whelan (2012) by allowing for cross-sectional and serial correlation in the forecast errors that can otherwise lead to incorrect inference. Taken together, the findings in our study and Hounyo and Lahiri (2023) contrast with previous evidence

³ We use the term "target variable" to denote the combination of an outcome variable (e.g., GDP growth) and forecast horizon (e.g., one-year-ahead horizon).

presented by Kenny, Kostka and Maserà (2014) and Meyler (2020) for the ECB-SPF and D’Agostino, McQuinn, and Whelan (2012) for the US-SPF.

There are, however, several important differences between our study and Hounyo and Lahiri (2023). Here we provide a short comparison and defer a more detailed discussion until Section II. Beyond methodology and data sets, one difference pertains to focus. The empirical framework of Hounyo and Lahiri (2023) is designed to test for equal predictive performance within a target variable. In contrast, we develop and conduct a more stringent test of comparable forecast behavior by considering both first and second moments of predictive performance. The ability to discriminate between participants who may display equal accuracy but differ along other dimensions of forecast behavior is an important extension in the evaluation of expectations models. Our investigation also extends beyond the inter-forecaster comparisons in Hounyo and Lahiri (2023) by considering intra-forecaster comparisons across target variables as an additional basis to explore the issue of the interchangeability of forecasters.

Another difference with Hounyo and Lahiri (2023) concerns the rank ordering of participants. While their approach allows for the identification and ranking of participants with superior or inferior forecasting skills during the sample period, it is silent on whether the ordering is stable over time. In contrast, our approach allows for a deeper investigation into the behavior of the rank orderings. We find that the bulk of the rank ordering of participants changes with variation in the forecasting environment. This result suggests that forecaster comparisons can be sensitive to the relative prevalence of tranquil and volatile episodes in a selected sample period and offers a cautionary note for studies that assume rank orderings are largely stable.

We conclude that models featuring information rigidities and their implication for forecaster interchangeability are not consistent with observed features of the ECB-SPF. The mean and variation over time in participants’ accuracy do not align with each other on average. In addition, we document systematic patterns in individual predictive performance that are strongly linked to the degree of difficulty in the forecasting environment. While such behavior could reflect heterogeneity in participants’ loss functions or the use of different models, a deeper exploration into this line of research is beyond the scope of this paper.⁴ We do, however, investigate the possibility of non-uniform processing capacity across forecasters that is, in turn, related to differential private information [Clements (2022)]. Taken together, our study principally contributes to a large literature

⁴ While strategic behavior could offer another explanation for these features, the anonymity of the ECB-SPF forecasters would likely rule out this explanation.

that uses survey data to inform the ongoing development of models of expectations formation, with particular focus on uncovering new facts about predictive performance.

The paper is organized as follows. The next section discusses the modeling strategy and estimation framework used for the empirical analysis. Section III provides a summary of the literature evaluating various aspects of professional forecasters' predictive performance. Section IV describes the ECB-SPF data. Section V reports the estimation results and documents systematic patterns in participants' relative predictive performance and how it varies with the difficulty of the forecasting environment. This section also explores the behavior of the rank ordering of forecasters. Section VI concludes by discussing the implications of our findings.

II. Modeling Strategy and Estimation Framework

Our modeling strategy and estimation framework are motivated by the survey-based predictive performance metric introduced by D'Agostino, McQuinn, and Whelan (2012) for the US-SPF and adopted by Hounyo and Lahiri (2023). A key aspect of both analyses concerns the challenge of evaluating predictive performance when there is time variation in the forecasting environment. Specifically, participants generating the same prediction error in different periods will not reflect equal predictive ability if forecasting in some periods is easier/more difficult as compared to others. The evaluation of predictive performance is further complicated in an unbalanced panel setting due to the entry and exit of participants either on an intermittent or permanent basis.

To account for both considerations, D'Agostino, McQuinn, and Whelan (2012) originally proposed an adjustment to conventional predictive performance metrics. Their approach begins by constructing a normalized forecast error statistic for each variable, period, and participant. While we discuss their methodology within the context of point forecasts, it can also be applied to density forecasts. Abstracting from details related to data and survey features, the normalized squared error statistic for participant j , $(\bar{e}_{t+h|t}^j)^2$, is given by:

$$(1) \quad (\bar{e}_{t+h|t}^j)^2 = \frac{(e_{t+h|t}^j)^2}{(1/N_t) \sum_{i=1}^{N_t} (e_{t+h|t}^i)^2} = \frac{(e_{t+h|t}^j)^2}{\overline{(e_{t+h|t})^2}}$$

where $e_{t+h|t}^j$ is participant j 's forecast error associated with the survey point prediction in period t and the realization of the target variable in period $t+h$, and $\overline{(e_{t+h|t})^2}$ is a measure of average forecast performance defined over the N_t survey participants in period t . Importantly, the metric in (1)

depends on participant j 's forecast performance relative to the other forecasters. Consequently, the normalization is designed to control for changes in the forecasting environment by generating, for a given value of $(e_{t+h|t}^j)^2$, a value of $(\bar{e}_{t+h|t}^j)^2$ that is lower (higher) when forecasters are collectively less (more) accurate compared to periods when they are more (less) accurate.

For each forecaster, we can calculate a score by taking an average of the normalized squared error statistics. Letting T^j denote the total number of surveys in which participant j appears and T denote the total number of surveys conducted, the score of participant j is defined as:

$$(2) \quad S^j = \left(\frac{1}{T^j} \right) \sum_{t=1}^T (\bar{e}_{t+h|t}^j)^2 ,$$

where $(\bar{e}_{t+h|t}^j)^2$ is set to zero if participant j did not respond to that survey. Because the performance score in (2) is calculated as an average, it can account for a participant entering or exiting a survey.

D'Agostino, McQuinn, and Whelan (2012) derive a historical distribution of forecast performance using the score in (2) and the associated rank ordering of all participants. A test for equal ability proceeds by randomly reshuffling and reassigning individual forecasts of a given variable for a particular survey. The same procedure is applied to each survey, resulting in a new sequence of forecasts for each participant that can be used to calculate an overall score from (2) and construct a rank ordering. The process is repeated many times to generate a large number of simulated distributions of forecaster performance, with the test for equal ability comparing the historical distribution of forecast performance to the simulated distributions. Under the null hypothesis of equal ability, the historical distribution of forecast performance should lie within selected percentiles of the simulated distribution that serve as confidence intervals. D'Agostino, McQuinn, and Whelan (2012) find little evidence that the best forecasters are significantly better than others, although there is a relatively small group of forecasters that perform very poorly.

While the approach in D'Agostino, McQuinn, and Whelan (2012) is attractive, Hounyo and Lahiri (2023) argue that the independent nature of the resampling method used to generate the simulated distribution is problematic for two reasons. First, common aggregate shocks can generate cross-sectional dependence across participants. Second, overlapping forecast horizons can generate time-series dependence across participants. As Hounyo and Lahiri (2023) note, D'Agostino, McQuinn, and Whelan (2012) independently resample observations from one forecaster to another within the same survey which ignores possible cross-sectional dependence in participants' forecast errors. They also note that D'Agostino, McQuinn, and Whelan (2012) independently resample

observations from one period to another which rules out possible serial correlation in a participant's forecast errors even though two of the four target variables involve overlapping forecast horizons.

To address these two concerns, Hounyo and Lahiri (2023) propose an alternative to the bootstrap procedure of D'Agostino, McQuinn, and Whelan (2012). Their method applies a wild bootstrap to the vector containing all the individual forecast errors at each point in time. Their approach accounts for any cross-sectional and serial correlation in participants' forecast errors while preserving the unbalanced nature of the panel.⁵ Importantly, the application of their testing procedure to forecasts of GDP growth and inflation from the US-SPF strongly rejects the null hypothesis of equal predictive ability, overturning the findings of D'Agostino, McQuinn, and Whelan (2012). In particular, the results indicate there are systematic differences in forecasters' ability that extend beyond the best forecasters and include forecasters across all percentiles of the distribution of predictive performance.

While the issue of equal predictive ability in Hounyo and Lahiri (2023) has relevance for our investigation, there is scope for further exploration into the properties of individual forecast behavior. For example, the implications of IR models for forecaster interchangeability apply equally to the mean and variation over time in accuracy and, therefore, would argue for also taking higher moments into consideration. In addition, the presence of systematic patterns in predictive performance raises questions about the source(s) for this feature of the data as well as the possible impact of these patterns on the rank ordering of forecasters. Hounyo and Lahiri (2023) rank forecasters based on average predictive performance, but it is not clear whether this ranking is stable over time.

On a more general level, another aspect of Hounyo and Lahiri (2023) to consider centers on their use of the metric in (1). Specifically, the normalized metric in (1) involves an asymmetric treatment of accuracy at the individual level versus the aggregate level. That is, a forecaster who makes a relatively large error when the average forecast error is small will incur a large penalty, whereas a forecaster who makes a relatively small error when the average forecast error is large will not benefit much.⁶

⁵ See Hounyo and Lahiri (2023) for a more detailed discussion.

⁶ For example, if an individual's forecast error is 0.5 percentage point and the average forecast error is 2 percentage points, then the individual's normalized forecast error decreases to 0.25 percentage point and there is a 'benefit' of 0.25 percentage point. On the other hand, if an individual's forecast error is 2 percentage points and the average forecast error is 0.5 percentage point, then the individual's normalized forecast error increases to 4 percentage points and there is a 'penalty' of 2 percentage points.

Drawing upon the previous discussion, we propose a heterogeneous panel data model to describe the predictive performance of survey participants. The empirical analysis is based on the following specification for the forecast performance of each participant j and survey in period t :

$$(3) \quad FP_{t+h|t}^j = \alpha_j + \lambda_j \left(\overline{FP}_{t+h|t} \right) + \varepsilon_{t+h|t}^j$$

where $FP_{t+h|t}^j$ and $\overline{FP}_{t+h|t}$ denote a forecast performance (FP) metric at the individual and cross-sectional average level, respectively, and $\varepsilon_{t+h|t}^j$ is a mean-zero error term. For the moment, we only note that lower (higher) values of the FP and \overline{FP} measures denote higher (lower) forecast accuracy and defer a more detailed discussion of the specific metrics until Section IV.

There is a close parallel between the specifications in (3) and (1) such that the panel data model can be viewed as a linear regression-based analogue to the normalized metric used in D’Agostino, McQuinn, and Whelan (2012) and Hounyo and Lahiri (2023). There are, however, several advantages to our empirical framework. First, the participant-specific intercept and slope allow for a deeper exploration into the nature of heterogeneity. Specifically, we can evaluate individual forecast performance through two channels: an individual fixed effect α -- which captures the component of a forecaster’s performance that is time-invariant -- and a time-varying component $\lambda(\overline{FP})$ -- which captures the component that depends on the degree of difficulty in the forecast environment. In addition, the linear specification in (3) does not maintain the asymmetric treatment of forecast accuracy in (1) at the individual level versus the aggregate level.

Our empirical framework also accounts for cross-sectional dependence in the data. Specifically, the inclusion of the cross-sectional average of predictive performance (\overline{FP}) in the individual regressions in (3) allows us to interpret our empirical framework within the context of the common-correlated effects (CCE) estimator of Pesaran (2006). As shown by Pesaran (2006), averaging the dependent variables in a panel data model at a point in time yields a proxy for an unobserved common component that can control for cross-sectional correlation across units.⁷ In the context of the ECB-SPF, the movements in \overline{FP} capture the effect of aggregate shocks that generate higher or lower accuracy across participants in a period and thereby also provide a very natural way to describe time variation in the difficulty of the forecasting environment.

⁷ See Pesaran (2006) for a more detailed discussion.

We calculate robust standard errors for the estimated parameters in (3) by applying the Newey-West (1987) covariance matrix modified for use in a panel setting to account for autocorrelation and conditional heteroscedasticity in the data. As previously discussed, the issue of time series dependence arises when the data involve overlapping forecast horizons which is relevant for our analysis.

The heterogeneity admitted by the specification in (3) also allows for a more detailed comparison of the statistical features of participants' accuracy and an evaluation of their alignment. Specifically, we can consider both first and second moments of accuracy as a basis to investigate the issue of comparable forecast behavior. As shown in the Appendix, the restriction $\alpha_j = 0$ and $\lambda_j = 1$ provides a test for distributional homogeneity which involves a joint test of equal predictive performance and equal variance of the predictive performance metric across participants.⁸ The testing procedure formalizes the idea that if forecasters are interchangeable, then over time their observed behavior should be indistinguishable from that of the consensus forecast.

In addition to the test for distributional homogeneity, another attractive feature of our empirical framework is that we can examine the estimated parameter pairings $(\hat{\alpha}_j, \hat{\lambda}_j)$ for evidence of other distinguishing patterns in participants' predictive performance within a target variable. A property of the regression equation (3) is that the estimated values for α and λ across participants will be centered around 0 and 1, respectively. As shown in Figure 1, we can partition the parameter space into four quadrants which affords an extremely intuitive way to visualize features of participants' predictive performance. If the estimated parameter pairings are not distributed randomly across the quadrants in Figure 1, then one possibility is that a scatterplot of the estimated parameter pairings principally run from the lower-left ($\alpha < 0, \lambda < 1$) quadrant up through the upper-right ($\alpha > 0, \lambda > 1$) quadrant. In this configuration, the lower-left quadrant would identify participants who are more accurate on average than their peers irrespective of the forecasting environment, while the upper-right quadrant would identify participants who are less accurate on average than their peers irrespective of the forecasting environment.⁹

A second possibility is that the estimated parameter pairings principally run from the upper-left ($\alpha < 0, \lambda > 1$) quadrant down through the lower-right ($\alpha > 0, \lambda < 1$) quadrant, implying that relative predictive performance varies with the forecasting environment. Specifically, participants in

⁸ We would like to thank an anonymous referee for bringing this point to our attention.

⁹ Recall that lower (higher) FP values are associated with higher (lower) individual forecast accuracy and lower (higher) FP values are associated with tranquil (volatile) forecasting environments.

the upper-left quadrant are relatively more accurate in a tranquil environment and then relatively less accurate as the environment becomes more volatile. The opposite holds for participants in the lower-right quadrant. Additionally, the quadrant scatterplot is informative about the dispersion of the estimated parameter pairings and their correspondence with the nature of the forecasting environment.

Estimation of (3) also allows us to generate a forecast performance profile for each participant based on the predicted values of the individual regressions (\overline{FP}) where:

$$(4) \quad \overline{FP}_{t+h|t}^j = \hat{\alpha}_j + \hat{\lambda}_j \left(\overline{FP}_{t+h|t} \right)$$

The behavior of the performance profiles depends on the forecasting environment and the quadrant location of the estimated parameter pairings. As we vary the value of \overline{FP} , the performance profile of participants located in the lower-left quadrant will run below the profile of those located in the upper-right quadrant without crossing. In the case of participants in the upper-left quadrant and lower-right quadrant, however, changes in the forecasting environment can cause their profiles to cross and generate variation in the rank orderings of predicted forecast accuracy. Our ability to observe movements in participants' rank orderings and gauge their stability over time provides a deeper insight into forecaster performance compared to other studies and highlights another important contribution of the analysis.

Our empirical framework also lends itself to making intrapersonal comparisons of the participants. Earlier discussion of IR models highlighted the implication that individuals should not display systematic patterns in their forecast behavior across target variables. To investigate this issue, we consider one approach that focuses on the quadrant location of a participant's estimated parameter pairing. Specifically, Section V describes a simulation exercise to assess if the quadrant locations for a participant's estimated parameter pairings are similar across target variables.

We also consider a second approach that examines the relationship between a participant's overall forecast performance relative to the consensus across target variables. Specifically, we construct the following metric for individual j for each target variable:

$$(5) \quad \left(\overline{FP}^j - \overline{FP} \right) = \left(1/T^j \right) \sum_{t=1}^T \left(FP_{t+h|t}^j - \overline{FP}_{t+h|t} \right)$$

where $FP_{t+h|t}^j$ is set to $\overline{FP_{t+h|t}}$ if participant j did not respond to that survey. The metric in (5) is similar to the score described in (2) and provides an assessment of a participant’s average relative forecast performance. That is, it indicates how a participant’s predictive performance compares to the cross-sectional average over time, with negative (positive) values associated with higher (lower) overall accuracy. The findings that a participant’s estimated parameter pairings tend to locate in the same quadrant and that the metric from (5) is correlated across target variables would indicate commonalities in a participant’s forecast behavior and offer evidence of deviations from IR models.

One last consideration is that the specification in (3) nests two alternative approaches previously used to capture the effects of aggregate shocks. Specifically, the normalization procedure proposed by D’Agostino, McQuinn, Whelan (2012) and maintained by Meyler (2020) and Hounyo and Lahiri (2023) corresponds to the restriction that the α_j ’s are jointly equal to zero, while the application of time fixed effects adopted by Kenny, Kostka, and Masera (2014) corresponds to the restriction that all of the λ_j ’s are equal. To preview our findings, the data reject both alternative approaches designed to control for variation in the forecasting environment.

Taken together, our modeling strategy provides a unified empirical framework to analyze the predictive performance of ECB-SPF participants and to inform models of the expectations formation process. Regarding IR models, our approach affords several avenues to analyze and characterize patterns in predictive performance and to determine if those patterns are consistent with interchangeable behavior on the part of participants. Moreover, our approach uses conventional estimation and testing procedures, as well as accounts for a range of econometric issues arising from the nature of the survey instrument and data. Importantly, changes in the predictability of a target variable do not present a challenge or require the adoption of some type of sub-sample analysis. Rather, time variation in the forecasting environment is an integral element in our methodology and plays a central role in our ability to compare and to contrast various features of participants’ predictive performance.

III. Literature Review

Our findings make several contributions to the existing literature on the expectations formation process. One area of interest focuses on the use of a panel data framework to explore different aspects of the predictive performance of ECB-SPF participants. Kenny, Kostka, and Masera (2014) compare the predictive performance of individual-level ECB-SPF density forecasts to density forecasts from a set of simple alternative benchmark models. Their results indicate

significant time variation in the forecast accuracy of participants relative to the benchmark models. Examining the link between the moments of density forecasts and density forecast performance, Kenny, Kostka, and Masera (2015a) find that forecast performance could be improved if participants corrected a downward bias in their reported variances. Kenny, Kostka, and Masera (2015b) report that predictive performance differs across forecasting tasks, with surveyed densities being much more informative about direction-of-change predictions than high and low outcome events. We also find significant time variation in the relative performance of both point and density forecasts. Moreover, we extend previous work on the determinants of differential predictive performance of ECB-SPF participants by identifying changes in the forecast environment as a new and important channel of influence that induces time-variation in the rank orderings of the panel.

Meyler (2020) applies the bootstrapping and Monte Carlo simulation techniques of D’Agostino, McQuinn, and Whelan (2012) to examine the issue of equal predictive performance for ECB-SPF point forecasts. He correctly notes that the testing procedure relies on participants’ forecast errors being uncorrelated across periods. When the data in (1) involve overlapping forecast horizons ($b > 1$), the conventional application of the testing procedure is not valid because of autocorrelation in the forecast errors. To remedy this situation, Meyler (2020) proposes separating the data across nonoverlapping forecast horizons. A drawback of this approach is that it restricts his analysis to SPF rounds that are four quarters apart which dramatically lowers the power of the testing procedure because of the reduced time series dimension of the data. An attractive feature of our approach is that it does not require the panel to be separated into sub-samples, thereby allowing us to exploit efficiency gains from using all the information on participants in a collective manner.

Another area of interest is heterogeneity in the forecast features of the ECB-SPF. Kenny, Kostka, and Masera (2014, 2015a, 2015b) find considerable heterogeneity in the performance of the surveyed densities, while Meyler (2020) finds little evidence of participants who perform significantly better or worse than their peers in terms of point forecasts. Our empirical framework provides an extremely flexible approach to investigate heterogeneity in predictive performance across multiple dimensions such as target variables and the quadrant location of a participant’s estimated $(\hat{\alpha}_j, \hat{\lambda}_j)$ parameter pairings. Moreover, we apply our empirical framework to both point and density forecasts as a robustness check and find that predictive performance displays stronger correlation patterns for the surveyed density forecasts than the point forecasts. Abstracting from other considerations, this result may help to explain some of the conflicting evidence reported in Kenny, Kostka, and Masera (2014, 2015a, 2015b) and Meyler (2020).

Our paper also contributes to a related literature that focuses more broadly on systematic patterns in survey-based forecasts. Bruine de Bruin et al. (2011) report strong evidence of persistence in individual participants' relative levels of uncertainty from the Federal Reserve Bank of New York Survey of Consumer Expectations. Boero, Smith, and Wallis (2015) examine the Bank of England Survey of External Forecasters and find significant persistence in the relative levels of point forecasts and uncertainty. Clements (2022) documents persistence in the relative levels of accuracy and disagreement of point forecasts for the US-SPF, while Rich and Tracy (2021) document persistence in the relative levels of disagreement and uncertainty for the ECB-SPF. While our study shares a similar motivation to Clements (2022), there are notable differences. For example, Clements (2022) restricts his analysis to point forecasts from the US-SPF and relies on a rank correlation test applied to two sub-samples to assess systematic patterns in individual forecast behavior.

Finally, our modeling strategy is closely related to the work of Qu, Timmermann, and Zhu (2019, 2021) that uses a panel data framework to analyze forecast accuracy. They consider various approaches to separate the importance of common shocks from idiosyncratic, individual-specific shocks. Our analysis differs in two important respects. The first is in terms of dimensions of the data. The modeling framework and testing procedures in Qu, Timmermann, and Zhu (2021) are designed for a large cross-section. However, our examination of the ECB-SPF only includes three variables – real GDP growth, inflation, and unemployment – which is too small for applying their methods. The second is in terms of focus. A key issue of interest in Qu, Timmermann, and Zhu (2019) is the identification of participants with superior forecasting skills. Consequently, their methodology involves the consideration of predictive performance across multiple dimensions and the assessment of a very large number of pairwise comparisons. In contrast, our interest is not in a detailed exploration aimed at an overall ranking of forecasters. Rather, the inter- and intrapersonal comparisons in our analysis are much more limited in scope and are more narrowly directed at assessing their consistency with expectations models that predict the absence of systematic patterns at the individual level.

IV. The European Central Bank Survey of Professional Forecasters

The ECB-SPF began in January 1999 and provides a quarterly survey of euro area forecasts. The survey draws its pool of panelists from both financial and nonfinancial institutions, with most, but not all, located in the euro area. Meyler (2020) notes that the principal aim of the survey is to solicit expectations about real GDP growth, inflation, and unemployment, although the questionnaire also contains a noncompulsory section asking participants for their expectations of

other variables and to provide qualitative comments that inform their quantitative forecasts.¹⁰ The ECB-SPF asks panelists for forecasts at short-, medium- and longer-term horizons, including both “rolling” and “calendar year” variants. The survey is fielded in January, April, July, and October, with approximately 55 panelists on average responding per survey. For additional details about the ECB-SPF, see Garcia (2003) and Bowles et al. (2007).

We examine forecasts for real GDP growth, HICP inflation, and the unemployment rate. This choice partly reflects the structure of the survey instrument that asks respondents to submit both point- and density-based forecasts for these three macroeconomic variables.¹¹ Because Kenny, Kostka, and Masera (2014, 2015a, 2015b) restrict their analyses to surveyed density forecasts and Meyler (2020) restricts his analysis to surveyed point forecasts, our inclusion of both types of forecasts offers an important robustness check. For the density forecasts, participants report their subjective probability distribution of forecasted outcomes as a histogram using a set of intervals provided in the survey. While the ECB-SPF occasionally changes the number of closed intervals for the histogram, it has essentially maintained a common bin width for the closed intervals throughout its history.¹²

Regarding forecast horizons, we examine point and density forecasts that involve rolling one-year-ahead and one-year/one-year-forward horizons. Compared to calendar year horizons, an advantage of the rolling horizons is that the horizon length remains constant through time and allows us to treat the data as quarterly observations on a set of individually homogeneous series. As Garcia (2003) notes, there is a temporal misalignment between the target variables because of differences in the data frequency and publication lags of the variables. Specifically, real GDP growth is published quarterly with a two-quarter lag, while HICP inflation and the unemployment rate are published monthly with a one-month and a two-month lag, respectively.¹³

Our study analyzes surveys conducted from 1999:Q1–2018:Q3, with forecast evaluation for all series ending in 2019:Q3. The ECB-SPF, like other surveys, has experienced entry and exit of

¹⁰ The additional expectations are for variables such as wage growth, the price of oil, and the exchange rate.

¹¹ The ECB-SPF is among a small but growing number of surveys that solicit both point and density forecasts. Other notable surveys include the US-SPF (published by the Federal Reserve Bank of Philadelphia), the Bank of England Survey of External Forecasters, and Federal Reserve Bank of New York Survey of Consumer Expectations.

¹² The only deviation in this design started with the 2020:Q2 survey in response to the COVID-19 outbreak. The (nearly) constant interval width of the ECB-SPF density forecasts contrasts with the US-SPF density forecasts, which have experienced periodic changes in interval widths.

¹³ For example, the 2010:Q1 survey questionnaire asks respondents to forecast one-year-ahead output growth from 2009:Q3–2010:Q3. For HICP inflation, the corresponding forecast horizon is December 2009–December 2010. For the unemployment rate, the corresponding forecast is for November 2010.

respondents over time. In addition, occasionally participants do not respond to a questionnaire or to individual items within the questionnaire. As noted by Meyler (2020), participants provide the highest number of forecasts for HICP inflation and the lowest number for unemployment, with the number of forecasts at the one-year-ahead horizon exceeding that at the one-year/one-year-forward horizon. Participants also report more point forecasts than density forecasts. Given the unbalanced panel structure of the ECB-SPF, we only include participants at each individual target variable/horizon who provide at least 50 forecasts.¹⁴ Further, we only consider matched point and density forecasts to maintain comparability across the types of forecasts. Consequently, the number of participants varies from 34 (HICP inflation at the one-year-ahead horizon) to 21 (unemployment rate at the one-year/one-year-forward horizon).

An important issue for the assessment of predictive performance is the choice of data vintage used to construct realizations of the target variables. As is the case for most macroeconomic data for most countries, euro-area macroeconomic statistics tend to be revised from preliminary releases. Consequently, a choice must be made about the relevant release associated with a participant’s forecast. Following Meyler (2020), we construct realizations of the target variables for HICP inflation and the unemployment rate using monthly data from the first full release.¹⁵ For real GDP growth, we construct realizations of the target variables using quarterly data from the second estimate. We have considered other approaches to construct realizations of the target variables as additional robustness checks.¹⁶

Another important issue for the assessment of predictive performance is the choice of point and density forecast accuracy measures. For the point forecasts, we adopt the absolute error as the metric:

$$(6) \quad \textit{POINT} FP_{t+h|t}^j = \left| X_{t+h} - E_t^j [X_{t+h}] \right|$$

where X_{t+h} denotes the realized value of the relevant ECB-SPF target variable in period $t+h$ and $E_t^j[X_{t+h}]$ denotes the reported point forecast from participant j in the survey at date t .

¹⁴ We have also experimented with a lower threshold of 40 participants and obtained similar results.

¹⁵ For example, if the target variable is one-year-ahead HICP inflation, we use the first full release reporting the value of the price index in month $t+12$. The same release is used to obtain the value of the price index in month t .

¹⁶ We used current vintage data as one robustness check. As another robustness check, we construct growth rates using the first full release to obtain the value of the price index in month t and the second estimate for the level of real GDP in quarter t . The results changed very little using these alternative approaches.

For the density-based accuracy measure, we adopt the absolute rank probability score (ARPS) as the metric:

$$(7) \quad \text{DENSITY } FP_{t+h|t}^j = \frac{1}{k_t - 1} \sum_{i=1}^{k_t} \left| \sum_{g=1}^i p_t^j - \sum_{g=1}^i I_{t+h} \right|$$

where we assume there are k_t bins associated with the histogram for the survey at date t , ${}_g p_t^j$ is the probability assigned by respondent j to the g^{th} bin, and ${}_g I_{t+h}$ denotes an indicator variable that takes a value of one if the actual outcome in period $t+h$ is in the g^{th} interval of the histogram from the survey at date t . The ARPS has the property that a participant receives “credit” by assigning probability in bins close to the bin containing the actual outcome.¹⁷

The evaluation of the ECB-SPF density forecasts requires additional discussion beyond the selections of data vintage and metrics. To the extent that respondents place any probability in either open interval, the manner chosen to close off the open intervals will affect the value of the forecast performance metric in (7). We follow a common—although ad hoc—assumption and close the exterior open intervals by assigning them twice the width of the interior closed intervals. We also need to address the issue of the location of probability mass associated with the density forecasts. We again draw upon common practices and assume that the probability mass is distributed uniformly within each bin of the histogram. Finally, we exclude the 2009:Q1 one-year-ahead real GDP growth density forecast data because many respondents placed significant probability in the lower open interval of the histogram in this survey.¹⁸

V. Empirical Results

We begin by examining the behavior of the (cross-sectional) average forecast performance metrics (\overline{FP}) to compare predictability across the target variables as well as to identify tranquil and volatile episodes. Figure 2 plots the movements of (\overline{FP}) for the point forecasts and density

¹⁷ The squared norm is used in Meyler (2020) for point forecasts and in Kenny, Kostka, and Masera (2014, 2015a) for density forecasts. Compared to the absolute value norm in (6) and (7), the squared norm is more sensitive to outliers and the manner used to close the exterior open intervals for density forecasts. For robustness, we also used the squared norm and found similar results.

¹⁸ For this survey, the significant probability mass at the lower open interval corresponded to a growth rate of “-1 percent or less” and was due to the survey design of the density forecasts and its inability to provide sufficient coverage for the pessimistic point predictions of output growth. For individuals who either reported point predictions below -1 percent or wanted to indicate significant downside risk, they assigned most of their probability to the open-ended interval. See Abel et al. (2016) for further discussion.

forecasts of real GDP growth, HICP inflation, and the unemployment rate at the one-year-ahead horizon, while Figure 3 provides the corresponding information at the one-year/one-year-forward horizon. The metrics are plotted based on the realization of the target variable, with gray bars indicating recessions as determined by the Euro Area Business Cycle Dating Committee of the Center for Economic Policy Research.¹⁹

As shown, there is generally a close correspondence between the point and density forecast performance metrics for the same target variable. While the difficulty of forecasting outcomes around the time of the global financial crisis and the euro-area debt crisis is evident, the data indicate other episodes associated with sizable forecast errors that are not uniform in their timing across the target variables. Consequently, there is sufficient variability in the forecasting environments to mitigate concerns that our results may be largely driven by just a few events.

As for the pattern of the forecast errors, they are highest for real GDP growth around the time of the global financial crisis. For HICP inflation, they are also highest around the time of the global financial crisis as well as elevated toward the beginning of the sample and during the middle of the last decade. For the unemployment rate, the forecast errors are again largest around the time of the global financial crisis, although they are also elevated at the beginning of the sample and around the time of the euro-area debt crisis.

Interpersonal Comparisons of Predictive Performance

We estimate the parameters in (3) using ordinary least squares (OLS), with standard errors computed using the Newey-West (1987) covariance matrix estimator modified for use in a panel data set.²⁰ Column 1 in Tables 1-2 presents formal tests for distributional homogeneity of the predictive performance metrics. Letting $\hat{\theta} = [(\hat{\alpha}_1, \hat{\lambda}_1), (\hat{\alpha}_2, \hat{\lambda}_2), \dots, (\hat{\alpha}_N, \hat{\lambda}_N)]$ denote the vector of estimated parameters of the model, we construct the following Wald test statistic for the joint null hypothesis that $\alpha_j = 0 \cap \lambda_j = 1$ for $j = 1, \dots, N$ participants in the panel for a specific target variable:

$$(8) \quad W = (\hat{\theta} - \theta_0)' [\text{var}(\hat{\theta})]^{-1} (\hat{\theta} - \theta_0)$$

¹⁹ For example, the metric associated with the forecasts of HICP inflation from 2015:Q1-2016:Q1 is plotted at 2016:Q1. While Figure 2 plots the value for the one-year-ahead point forecasts of real GDP growth in 2009:Q1, recall that the analysis does not include these data due to the exclusion of the matched density forecasts. Unlike the absolute error metric, the ARPS metric is restricted to fall in the range between 0 and 1.

²⁰ We allow the error terms to follow a fourth-order moving average process to account for the overlap of forecast horizons.

The values of the test statistic indicate strong evidence of systematic differences in the mean and variance of participants' forecast accuracy as we reject the null hypothesis at the 1 percent significance level in all cases except for the point forecasts of inflation at the one-year horizon.²¹

The test of distributional homogeneity involves a joint test of equal predictive performance and equal variance of the predictive performance metric. Because almost all investigations into forecast performance have focused exclusively on equal predictive ability, it would be of interest to explore how this issue may bear upon the reported rejections. Even if participants display equal predictive ability, our more restrictive testing procedure could lead to a rejection of the null hypothesis due to heterogeneity in the variance of the performance metric. To investigate this possibility, we apply the testing procedure of Hounyo and Lahiri (2023) to the ECB-SPF data and report the findings in the Appendix.

As shown, the evidence strongly rejects the null hypothesis of equal predictive ability at conventional significance levels.²² Consequently, the rejection of distributional homogeneity does not mask equal predictive ability among forecasters. It is also interesting to note that the pattern of rejections of equal predictive ability is remarkably similar to that documented by Hounyo and Lahiri (2023) for the US-SPF where the best forecasters as well as forecasters across the other percentiles of the distribution of predictive performance are more accurate than what would be expected by random chance using the bootstrap procedure.

Because our empirical framework nests two common approaches to control for variability in the forecasting environment, we also construct Wald test statistics for the validity of the normalized predictive performance metric ($\alpha_1 = \alpha_2 = \dots = \alpha_N = 0$) and for the use of time fixed effects ($\lambda_1 = \lambda_2 = \dots = \lambda_N$). Except for the point forecasts of inflation at the one-year horizon, the results in column 2 and column 3 strongly reject the normalized predictive performance metric and the use of time fixed effects to control for the effects of aggregate shocks, respectively. These findings offer an additional reason why our evidence of systematic differences in the predictive performance of ECB participants contrasts with the analyses of Kenny, Kostka, and Masera (2014) and Meyler (2020). Specifically, Kenny, Kostka and Masera (2014) use time fixed effects to control for changes in the forecasting environment., while Meyler (2020) adopts the normalization procedure.²³

²¹ The different degrees of freedom reflect the varying number of respondents meeting the participation restriction for the various target variables.

²² Following Hounyo and Lahiri (2023), we also consider the case of excluding forecasters who scored worse than the 80th percentile and found similar results. These findings are also reported in the Appendix.

²³ As previously discussed, the more general nature and greater parameter flexibility of our empirical framework allows the model estimates to account for changing relative performance rankings. The inability of

To gain insight into the relative accuracy of participants, Figures 4-5 display scatterplots and correlation coefficients (r) for the individual estimated parameter pairings $(\hat{\alpha}_j, \hat{\lambda}_j)$ for the point forecasts and density forecasts, respectively. The patterns are striking in their similarity across target variables. Because few of the estimated parameter pairings fall in the lower-left ($\alpha < 0, \lambda < 1$) and upper-right ($\alpha > 0, \lambda > 1$) quadrants, the visual evidence does not support the view that the ECB-SPF panel is comprised of participants who remain relatively more accurate and other participants who remain relatively less accurate across all forecasting environments. Instead, evidence of the estimated parameter pairings largely falling in the upper-left ($\alpha < 0, \lambda > 1$) and lower-right ($\alpha > 0, \lambda < 1$) quadrants indicates that participants' relative accuracy varies with the forecasting environment. Moreover, the patterns do not suggest clustering or that the negative relationship reflects the behavior of a few participants. Rather, the observations are dispersed within each of the two quadrants and display a comparable count across the two quadrants. We also observe the correlations are larger in absolute value for the density forecasts compared to the point forecasts.

Given the evidence documenting a strong link between a participant's predictive performance and the difficulty of the forecasting environment, it is natural to ask what might be driving this result. Here we consider one possible explanation that draws upon Clements (2022) and can be easily incorporated within our empirical framework. Specifically, Clements (2022) examines the US-SPF and investigates whether systematic differences in forecast accuracy are related to systematic differences between forecasters in their degree of contrarianism. Such would be the case if some forecasters receive superior private information, resulting in predictions that display greater contrarianism but also higher accuracy.

We investigate the relationship between forecast accuracy and contrarianism by pairing a participant's average relative forecast performance metric in (5) with an analogue for disagreement. Specifically, we construct the following measure of average relative disagreement:

$$(9) \quad \left(\overline{D^j - \bar{D}} \right) = \left(1/T^j \right) \sum_{t=1}^T \left(D_{t+h|t}^j - \overline{D_{t+h|t}} \right)$$

where

$$(10) \quad D_{t+h|t}^j = \left| E_t^j [X_{t+h}] - \overline{E_t [X_{t+h}]} \right|$$

the specifications in Kenny, Kostka and Masera (2014) and Meyler (2020) to capture this feature of the data may be another factor explaining why our conclusions differ.

and where $D_{t+h|t}^j$ is set to $\overline{D_{t+h|t}}$ if participant j did not respond to that survey. The metric in (9) indicates how a participant's disagreement compares to average disagreement across the responding survey group over time, where individual disagreement is described in (10) and is measured by the absolute value of the deviation between a participant's forecast and the consensus forecast.²⁴ Because a positive (negative) value in (9) reflects higher (lower) relative disagreement, we would expect a negative association with average relative forecast performance if superior information is the source for the persistent performance heterogeneity among forecasters.

Figure 6 displays scatterplots and correlation coefficients (r) for the individual pairings of disagreement and predictive performance for the point forecasts. As shown, there is a positive relationship of varying strength between disagreement and accuracy across target variables which suggests that more contrarian forecasters on average make less accurate forecasts. This finding is consistent with evidence in Clements (2022) for the US-SPF and does not support the conjecture that heterogeneity in forecast accuracy reflects some participants benefiting from superior private information.

Time-variation in Forecast Performance Profiles

An attractive feature of our empirical framework is that we can examine the implications of the estimation results for the forecast performance profiles of participants. To illustrate, we will initially select a participant from each of the four quadrants associated with a target variable. While any scatterplot can be used for the exercise, we select the one-year-ahead point forecasts of GDP growth because it is likely to be of particular interest.²⁵ Figure 7 depicts the data and the estimated regression line for each participant identified by the color-coded circles for the one-year-ahead GDP growth rate in the upper-left panel in Figure 4. As shown, the estimated regression lines display a very high fit to the data and suggest little reason to depart from the linear specification in (3).²⁶

Figure 8 plots the predicted forecast performance profiles of the same four participants and provides a visual investigation into their behavior as well as the incidence and nature of crossings that bear upon the issue of the stability of rank orderings. Using the estimated parameter pairing for

²⁴ Similar to the inclusion of \overline{FP} in (5), the inclusion of average disagreement ($\overline{D_{t+h|t}}$) in (9) is consistent with Clements (2022) who cites the importance of controlling for variation in the extent of disagreement over time.

²⁵ While the one-year-ahead point forecasts of inflation would also be of interest, recall that we do not reject the property of distributional homogeneity for this series.

²⁶ There is an outlier observation for three of the four participants associated either with realized GDP growth in 2008:Q3 or 2008:Q4. We exclude the relevant observation from the scatterplots (but not the reported R^2 values) in Figure 7 to enhance presentation of the data. The scatterplots and regression lines including all observations are provided in the Appendix.

each participant, we vary the average forecast performance metric (\overline{FP}) to trace out the performance profiles. Figure 8 also includes a 45-degree line indicating where individual predicted forecast performance (\overline{FP}) equals the cross-sectional average of forecast performance (\overline{FP}).

The resulting performance profiles closely align with our expected behaviors.²⁷ If participants were principally located in the lower-left (purple) and the upper-right (blue) quadrants, then this configuration would produce relatively stable rankings over time. This is shown by the purple and blue lines not crossing, with increases in the difficulty of the forecasting environment only acting to widen the gap between them. Because participants in the lower-left quadrant are predicted to be systematically more accurate than the average, the purple line always lies below the 45-degree line. The opposite holds for the participant in the upper-right quadrant. It is important to note that our illustration does not claim that the performance profiles cannot display crossings, but it does indicate that the crossings can only occur among participants located in the same quadrant.

As shown in Figures 4 and 5, most participants are located in the upper-left (yellow) and lower-right (red) quadrants. In contrast to the previous configuration, this configuration will produce highly variable rankings over time as performance profiles will display crossings beyond those involving participants located in the same quadrant. This is illustrated on a general level by the yellow and red lines crossing the 45-degree line which indicates a switch in the forecast accuracy of the participants relative to the cross-sectional average. Focusing on our selected participants, we see the yellow and red lines cross at 0.83 (the 59th percentile of \overline{FP}) as well as crossings with the participant from the lower-left quadrant at 0.41 (29th percentile) and 1.45 (90th percentile), respectively.²⁸ As shown, these crossings are associated with changes in rankings of participants as the forecasting environment evolves from low difficulty to extreme difficulty.²⁹

²⁷ Similar to Figure 7, there is a corresponding outlier value of \overline{FP} that we elect to exclude from the plots in Figure 8 to enhance presentation of the performance profiles. The performance profiles values using the full range of \overline{FP} values are also provided in the Appendix.

²⁸ While participants in the lower-left quadrant display forecasts that are systematically more accurate than the average, this does not imply that their forecasts always outperform those of individuals in other quadrants. Consequently, there is no inconsistency with the figure displaying the crossings of the performance profiles by the participants in the upper-left and lower-right quadrants. A similar point holds for participants in the upper-right quadrant.

²⁹ As shown in the corresponding figure in the Appendix, there is an additional crossing of the yellow and blue lines near the upper range of the \overline{FP} values. While a crossing point can always be calculated between the 45-degree line and the performance profile of a participant in the upper-left or lower-right quadrant, this may occur outside the range of \overline{FP} values in the sample.

To gain a better appreciation of the extent of this variability, we now consider all 33 participants associated with the one-year-ahead point forecasts of GDP growth. We construct rank orderings based on participants' forecast accuracy evaluated at eight values spanning the range of \overline{FP} values for this target variable. Table 3 reports the results, where the first column lists the forecaster IDs and the remaining columns moving from left to right indicate the rank ordering of each forecaster as the forecasting environment becomes more difficult. For ease of comparison, the order of the ID numbers is based on the initial ranking of forecasters at the lowest value of $\overline{FP} = 0.25$.

As shown in Table 3, the pattern of the rank orderings is consistent with the evidence from the scatterplot of the estimated parameter pairings in the upper-left panel of Figure 4. Some respondents largely maintain similar rankings either because they tend to be highly accurate (#22), highly inaccurate (#36), or close to the cross-sectional average most of the time (#54, #16). For most respondents, however, their rank orderings vary over the forecast environment. While this variation can reflect dramatic improvements (#94, #52) or dramatic declines (#37, #39) in predictive performance, it is more typical to observe individuals who become relatively more accurate (#24, #2) or relatively less accurate (#54, #98) as the forecasting environment turns more challenging.

There are two key takeaways that emerge from Figure 8 and Table 3. The first is that forecaster evaluations and comparisons may not be invariant to the relative prevalence of tranquil and volatile episodes in a selected sample period. The second is that the evidence may provide one explanation for the finding that it is difficult, *ex ante*, to devise forecast combination methods that beat a simple average.³⁰ While our analysis links predictive performance to variation in the forecasting environment, this feature may not be easily exploitable because of the inherent difficulty of predicting tranquil/volatile episodes in real time. If, however, we were to allow for the availability of some information on an *ex post* basis, then the Appendix describes a “performance weighting” combination scheme that consistently and significantly outperforms the equally-weighted consensus forecast.³¹

Intrapersonal Comparisons of Predictive Performance

The analysis up to this point has examined forecast data for the target variables in isolation. However, we can also investigate if there are commonalities in individual predictive performance across target variables. While the scatterplots in Figure 4 and Figure 5 show that the estimated parameter pairings principally lie in the $(\alpha < 0, \lambda > 1)$ and $(\alpha > 0, \lambda < 1)$ quadrants, they do not

³⁰ See Timmermann (2006) and Genre et al. (2013).

³¹ We would like to thank an anonymous referee for suggesting this exercise.

indicate the extent to which the pairings for a participant tend to locate in the same quadrant across target variables. Another consideration is the extent to which a participant's predictive performance for a target variable correlates with performance for other target variables. Previous discussion has noted the relevance of these additional considerations for IR models.

Our investigation into the location of the parameter pairings for a participant requires more than simply focusing on the quadrant associated with the estimates. We must also account for the uncertainty associated with the estimates. Consequently, we adopt Monte Carlo simulation techniques and generate 1,000 draws of the parameter pairings vector for each target variable and horizon using the estimated joint normal distribution for $\hat{\theta} = [(\hat{\alpha}_1, \hat{\lambda}_1), (\hat{\alpha}_2, \hat{\lambda}_2), \dots, (\hat{\alpha}_N, \hat{\lambda}_N)]$. We can then use the simulated distributions to calculate the percentage of simulated parameter pairings located in each quadrant for a participant.

Figure 9 and Figure 10 plot the distributions for the point forecasts and density forecasts, respectively, where the estimated pairings are color-coded in black and the simulated pairings are in gray. As shown, the distributions for the density forecasts are much tighter compared to the point forecasts.³² While we only make note of this difference at present, it would be interesting for future research to explore the reasons for this feature of the data. For example, it is possible that respondents make less use of rounding and report less judgmental density forecasts compared to point forecasts.³³ Relatedly, Glas and Hartmann (2022) analyze the US-SPF and ECB-SPF and note that survey participants who report rounded point forecasts differ from respondents who round probabilities for density forecasts.

Figure 11 and Figure 12 focus on the quadrant location of the parameter pairings associated with participants' point forecasts and density forecasts, respectively. The values report the highest fraction of simulated parameter pairings for a participant that fall in the same quadrant based on the 6,000 simulations (1,000 simulations for each of the six target variables). For purposes of comparison, we present the results for the 23 participants included in all six combinations of target variables.³⁴ Overall, the evidence in Figure 11 provides general support for the idea that a participant's parameter pairings tend to locate in the same quadrant as almost all the histogram bars exceed 40 percent. Using 50 percent as an arbitrary threshold, the histogram bars show that about a third of the participants exceed the threshold criterion.

³² The difference in precision may explain why rejections of various hypotheses in Tables 1-2 are stronger for the density forecast data.

³³ We would like to thank an anonymous referee for bringing this point to our attention.

³⁴ We have also extended the analysis to include the other participants in our study and the results are similar.

A very different picture emerges when we look at the density forecasts. There are now 19 participants who exceed the threshold criterion, with the calculated percentages notably higher than the 50 percent value in many cases. Particularly noteworthy are the two participants whose forecast behavior suggests their parameter pairings would almost always fall in the same quadrant across the six combinations of target variables. Compared to the point forecasts, the density forecasts indicate considerably more overlap in quadrant locations which is consistent with the evidence from Figure 9 and Figure 10 and provides another example of the different conclusions that can be drawn between the point and density forecasts.

We can also use the average relative forecast performance metrics in (5) to make various comparisons across the forecast data, where lower values again indicate better predictive performance. Because of the large number of comparisons, we only provide a summary of the results. Overall, we find that forecast performance correlates positively across horizons and outcome variables in almost all cases. There are, however, differences across some dimensions that are worth noting. One difference is that the density forecast data generate a much stronger association than the comparable point forecast data. The top panel in Figure 13 is representative of this finding and shows scatterplots of the average relative forecast performance metrics for inflation at the two forecast horizons for the point forecast data and density forecast data, respectively. While the point forecast data indicate a modest correlation of 0.43, the density forecast data indicate a correlation of 0.77 which is nearly twice as high. Looking across all pairwise combinations of target variables, the correlations for the point forecast data are typically in the 0.2-0.4 range, while the correlations for the density forecast data are in the 0.7-0.8 range.

Another feature of predictive performance that emerges is that the correlations are generally higher for the same target variable at different horizons than for different target variables at the same horizon. An ordered ranking of the correlations indicates that the lowest three values are associated with inflation and GDP growth at the two forecast horizons and GDP growth and unemployment at the one-year/one-year-forward horizon. In contrast, the highest three values are associated with unemployment (using both types of forecast data) and inflation at the two forecast horizons.

A further examination of forecast performance across the target variables reveals two other features. First, there tends to be a stronger correlation at the shorter horizon. The middle panel of Figure 13 shows scatterplots of the average relative forecast performance metrics for GDP growth and unemployment. Unlike the pattern at the one-year-ahead horizon, there is much less of a translation of forecast performance from unemployment into GDP growth at the one-year/one-

year-forward horizon. Second, there is less of a linkage between forecast performance for GDP growth and inflation than there is for GDP growth and unemployment. The bottom panel of Figure 13 shows the scatterplots of the corresponding average relative forecast performance metrics at the one-year-ahead horizon, where we again include the GDP growth/unemployment scatterplot to facilitate the comparison. For inflation and GDP growth, predictive performance shows a slightly negative relationship.³⁵ In the case of GDP growth and unemployment, however, there is a sufficiently meaningful positive association.

VI. Conclusion

This paper adopts the common correlated effects (CCE) estimator of Pesaran (2006) to investigate whether ECB-SPF participants can be viewed as interchangeable. While the behavior of professional forecasters is of interest by itself, our study draws further motivation from IR models and their implication that systematic patterns should not be evident in the forecast data. In addition to making comparisons of predictive performance across participants, we investigate the correlation patterns for an individual's predictive performance across parameter configurations and target variables. As a robustness check, we also consider the evidence from point forecasts and density forecasts.

Based on forecasts for output, inflation, and unemployment, we find strong evidence of systematic patterns in participants' predictive performance. Moreover, the patterns are not a consequence of differential innate ability, but instead are episodic in nature and directly linked to changes in the forecasting environment. By way of a simple narrative, our interpersonal analysis of predictive performance suggests that participants largely divide into two "camps": those who display relatively more accurate forecasts in low-variance times and those who do so in high-variance times. Consistent with this view, we find the rank orderings of participants shift over time and display considerable variability.

Our intrapersonal analysis of predictive performance indicates that the influence of the forecasting environment carries over to other features of the forecast profile of participants. Specifically, we find there are commonalities across parameter quadrants and target variables, with the density forecast data revealing greater similarities in individual forecast behavior. In terms of the narrative introduced above, participants tend to locate in the same "camp" which indicates that a

³⁵ This is the only instance where the average relative forecast performance metrics display a negative relationship.

participant's relative accuracy in a forecasting environment is positively correlated across target variables.

Overall, we conclude that the predictive performance of ECB-SPF participants reflects distinguishing behaviors that are inconsistent with the implications of IR models for interchangeability. The strong evidence of systematic patterns in predictive performance and their relationship to the nature of the forecasting environment is a new finding and reflects the capabilities and advantages of the CCE empirical framework.

It would be interesting and important to determine if these same empirical features are present in other long running panel survey data.³⁶ Our findings support further development of expectations models that can generate systematic patterns in key features of forecasters' behavior as well as account for the differential effects of the forecast environment on predictive performance. The opportunity to explore and identify the key underpinnings during such a development would serve as fertile ground for future research.

³⁶ The US-SPF would seem to be a natural candidate. It is unclear, however, if a parallel analysis can be conducted for the US-SPF because of differences in the survey instrument. Specifically, the US-SPF does not feature "rolling" forecast horizons. Such investigations, however, would not need to be restricted to professional forecasters and should be considered for surveys more generally.

Table 1			
Comparison of Predictive Performance Behavior of ECB-SPF Participants			
Point Forecasts			
$^{POINT}FP_{t+h t}^j = \alpha_j + \lambda_j \left(\overline{FP_{t+h t}} \right) + \varepsilon_{t+h t}^j$			
	Distributional Homogeneity	Normalization Approach	Time Fixed Effects
Point Forecast Data	$H_0 : \alpha_j = 0 \cap \lambda_j = 1 \forall j$	$H_0 : \alpha_1 = \alpha_2 = \dots = \alpha_N = 0$	$H_0 : \lambda_1 = \lambda_2 = \dots = \lambda_N$
GDP growth: one-year-ahead	$\chi^2(62) = 332.8^{**}$	$\chi^2(31) = 76.0^{**}$	$\chi^2(30) = 245.9^{**}$
GDP growth: one-year/one-year-forward	$\chi^2(58) = 244.7^{**}$	$\chi^2(29) = 68.1^{**}$	$\chi^2(28) = 143.9^{**}$
Inflation: one-year-ahead	$\chi^2(68) = 77.7$	$\chi^2(34) = 41.1$	$\chi^2(33) = 29.3$
Inflation: one-year/one-year-forward	$\chi^2(62) = 156.2^{**}$	$\chi^2(31) = 90.9^{**}$	$\chi^2(30) = 80.4^{**}$
Unemployment: one-year-ahead	$\chi^2(56) = 134.5^{**}$	$\chi^2(28) = 58.0^{**}$	$\chi^2(27) = 58.1^{**}$
Unemployment: one-year/one-year-forward	$\chi^2(48) = 225.8^{**}$	$\chi^2(24) = 42.0^*$	$\chi^2(23) = 100.1^{**}$

Note: Model parameters are estimated using ordinary least squares (OLS), with standard errors computed using the Newey-West (1987) covariance matrix estimator modified for use in a panel data set. The error terms to follow a fourth-order moving average process to account for the overlap of forecast horizons. Degrees of freedom are reported in parentheses.

** Significant at the 1% level

* Significant at the 5% level

Table 2			
Comparison of Predictive Performance Behavior of ECB-SPF Participants			
Density Forecasts			
$DENSITY_{FP_{t+h t}}^j = \alpha_j + \lambda_j \left(\overline{FP_{t+h t}} \right) + \varepsilon_{t+h t}^j$			
Density Forecast Data	Distributional Homogeneity $H_0 : \alpha_j = 0 \cap \lambda_j = 1 \forall j$	Normalization Approach $H_0 : \alpha_1 = \alpha_2 = \dots = \alpha_N = 0$	Time Fixed Effects $H_0 : \lambda_1 = \lambda_2 = \dots = \lambda_N$
GDP growth: one-year-ahead	$\chi^2(62) = 348.5^{**}$	$\chi^2(31) = 240.1^{**}$	$\chi^2(30) = 191.3^{**}$
GDP growth: one-year/one-year-forward	$\chi^2(58) = 371.2^{**}$	$\chi^2(29) = 173.1^{**}$	$\chi^2(28) = 100.6^{**}$
Inflation: one-year-ahead	$\chi^2(68) = 352.6^{**}$	$\chi^2(34) = 208.8^{**}$	$\chi^2(33) = 103.2^{**}$
Inflation: one-year/one-year-forward	$\chi^2(62) = 345.6^{**}$	$\chi^2(31) = 273.5^{**}$	$\chi^2(30) = 134.3^{**}$
Unemployment: one-year-ahead	$\chi^2(56) = 324.2^{**}$	$\chi^2(28) = 141.1^{**}$	$\chi^2(27) = 106.0^{**}$
Unemployment: one-year/one-year-forward	$\chi^2(48) = 138.6^{**}$	$\chi^2(24) = 69.8^{**}$	$\chi^2(23) = 74.6^{**}$

Note: Model parameters are estimated using ordinary least squares (OLS), with standard errors computed using the Newey-West (1987) covariance matrix estimator modified for use in a panel data set. The error terms to follow a fourth-order moving average process to account for the overlap of forecast horizons. Degrees of freedom are reported in parentheses.

** Significant at the 1% level

* Significant at the 5% level

Table 3

Rank Orderings of Forecast Accuracy: Point Forecasts of One-Year-Ahead GDP Growth

Forecaster ID	$\overline{FP} = 0.25$	$\overline{FP} = 0.50$	$\overline{FP} = 0.75$	$\overline{FP} = 1.00$	$\overline{FP} = 1.50$	$\overline{FP} = 2.00$	$\overline{FP} = 4.00$	$\overline{FP} = 6.00$
37	1	4	4	8	14	21	25	25
39	2	3	3	5	11	17	21	24
42	3	2	2	2	5	8	10	13
22	4	1	1	1	1	2	5	6
61	5	5	8	13	16	20	20	20
95	6	6	7	12	13	15	18	18
4	7	7	10	14	17	19	19	19
5	8	10	18	21	24	26	27	27
88	9	8	6	9	10	10	9	9
23	10	9	5	4	7	7	7	7
89	11	19	25	26	27	27	28	29
54	12	14	17	18	20	18	17	17
33	13	27	29	30	31	33	33	33
38	14	11	11	11	9	9	8	8
15	15	26	28	29	30	30	31	31
16	16	16	16	17	15	12	14	14
47	17	15	15	16	12	11	12	11
98	18	20	22	23	22	23	23	23
31	19	21	23	24	23	24	22	22
56	20	25	27	27	28	29	29	28
24	21	17	19	19	19	16	15	16
85	22	13	12	10	6	5	6	5
93	23	24	26	25	25	25	24	21
26	24	12	9	3	4	4	3	3
29	25	30	30	32	32	32	32	32
20	26	23	21	20	18	13	13	12
96	27	18	13	6	3	3	2	2
1	28	29	24	22	21	14	11	10
94	29	22	14	7	2	1	1	1
52	30	28	20	15	8	6	4	4
2	31	31	32	31	29	28	26	26
36	32	33	33	33	33	31	30	30
90	33	32	31	28	26	22	16	15

Figure 1

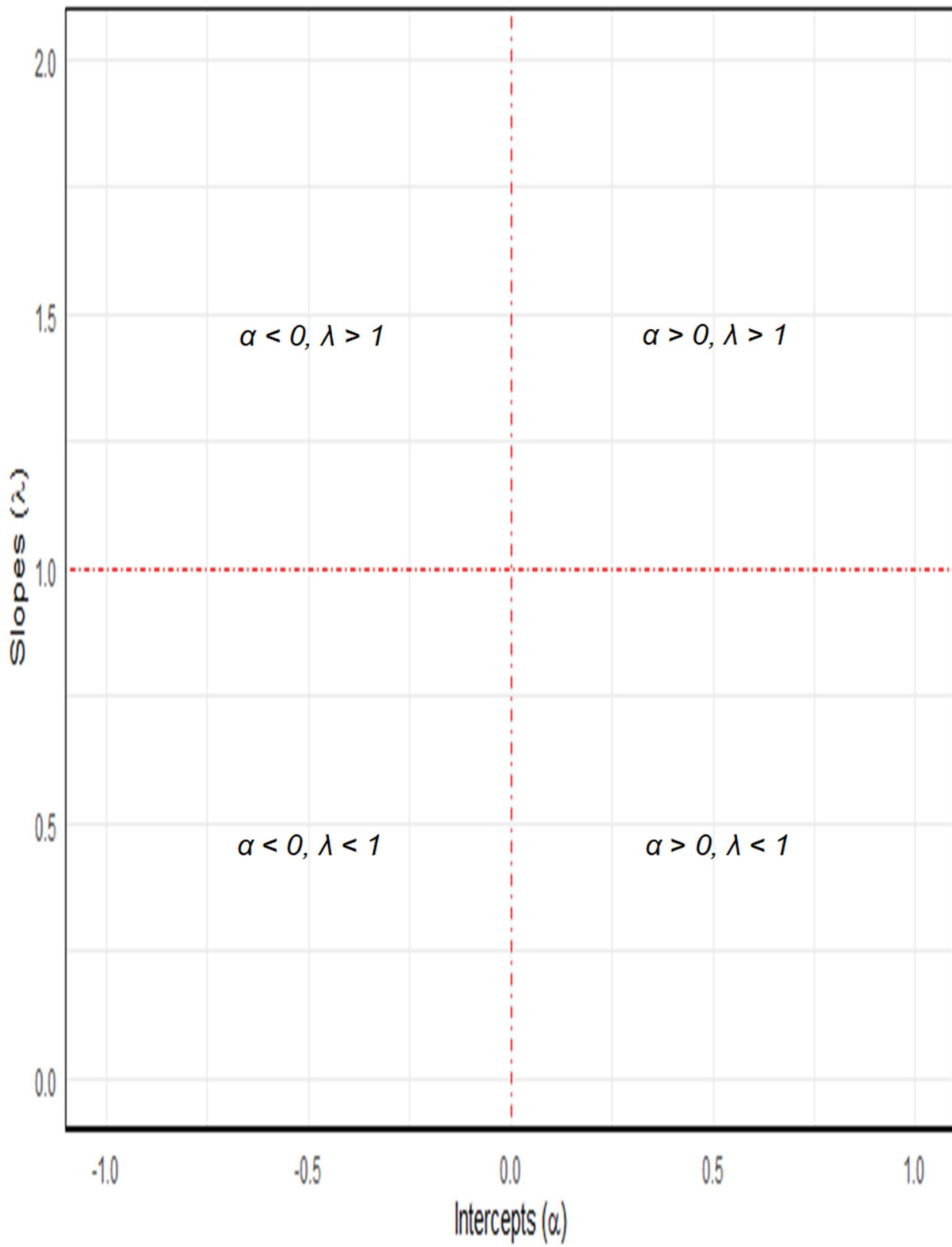


Figure 2

Average Forecast Performance: One-Year-Ahead Forecasts

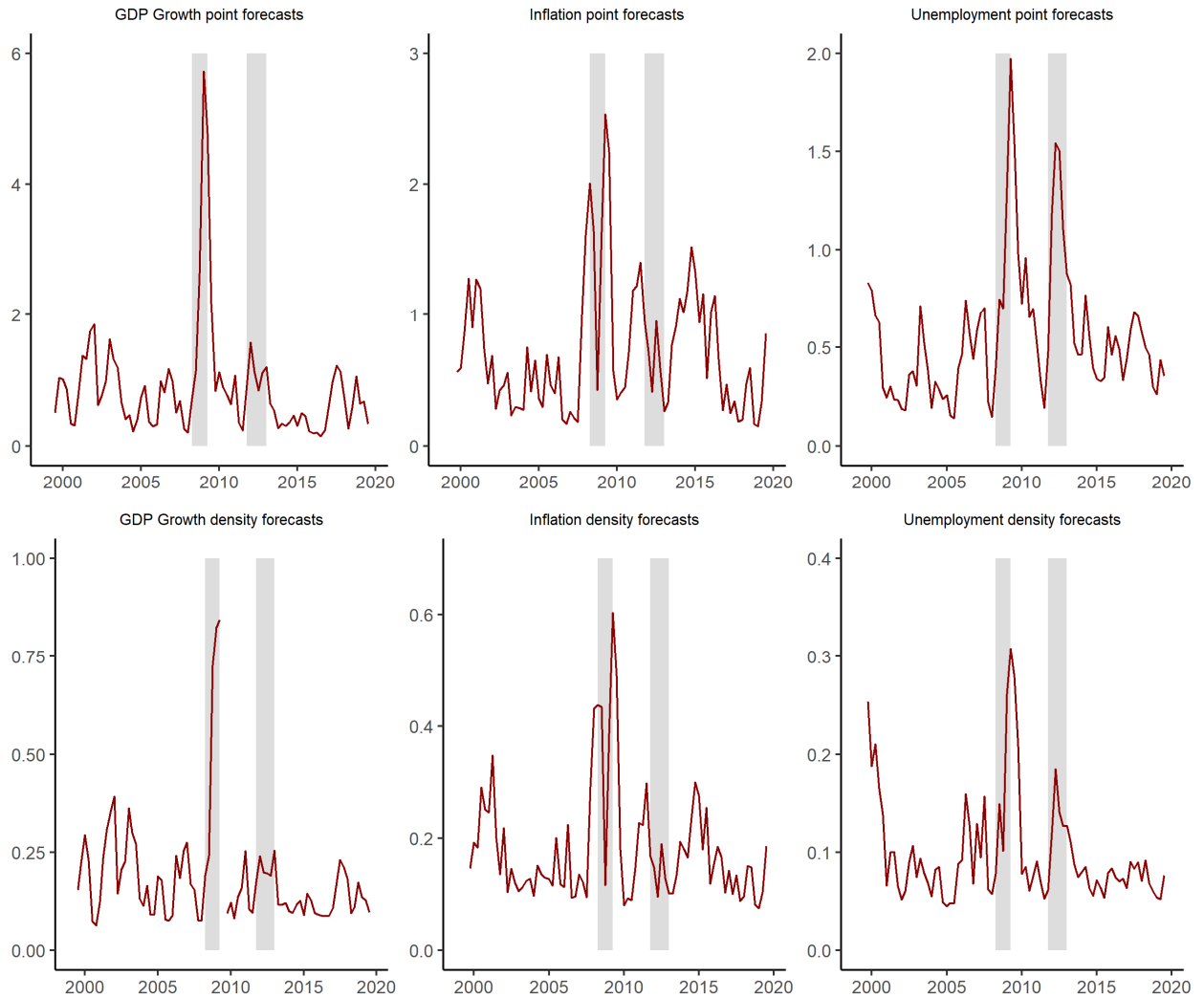


Figure 3

Average Forecast Performance: One-Year/One-Year Forward Forecasts

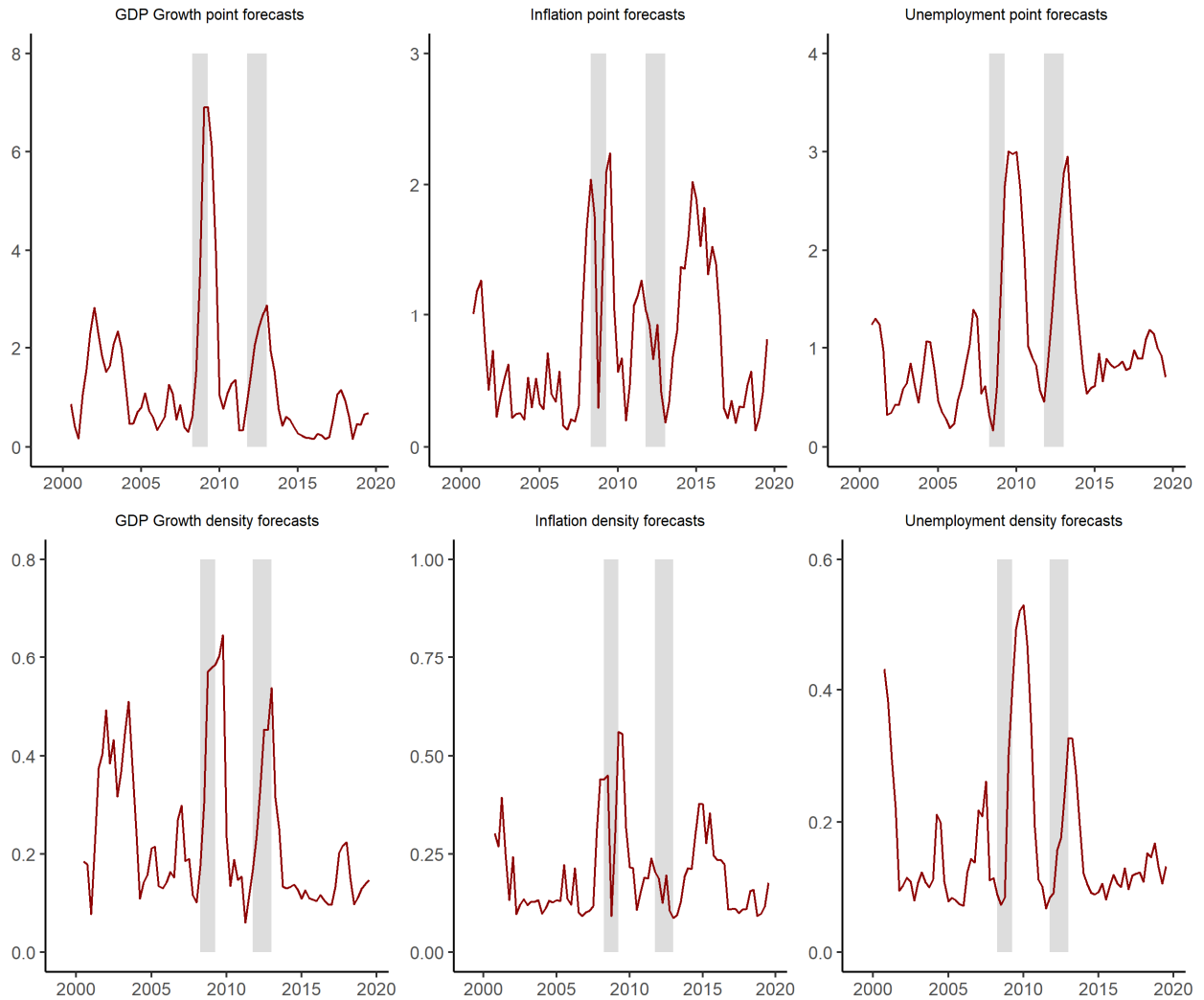


Figure 4

Estimated Parameter Pairings: Point Forecasts

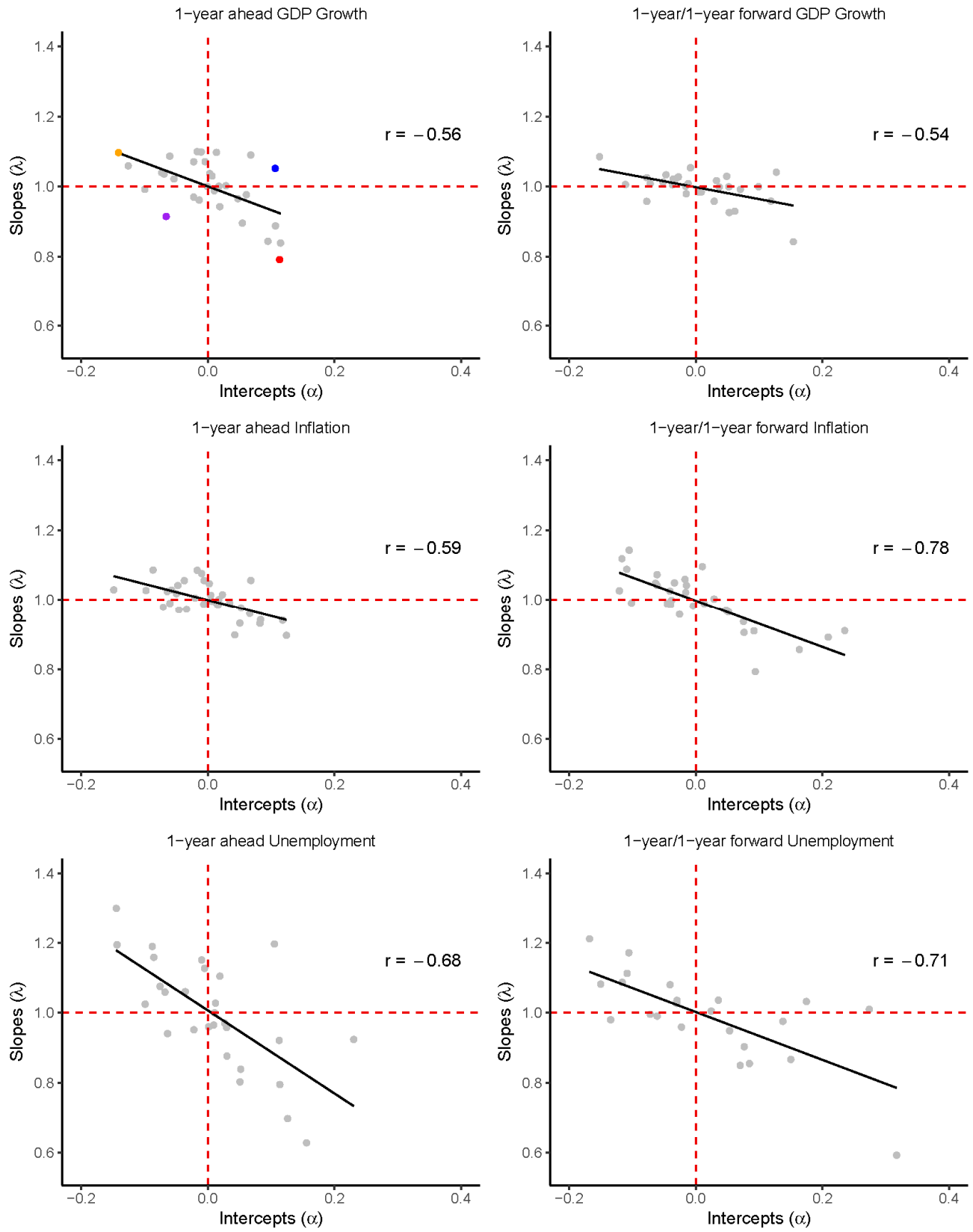


Figure 5

Estimated Parameter Pairings: Density Forecasts

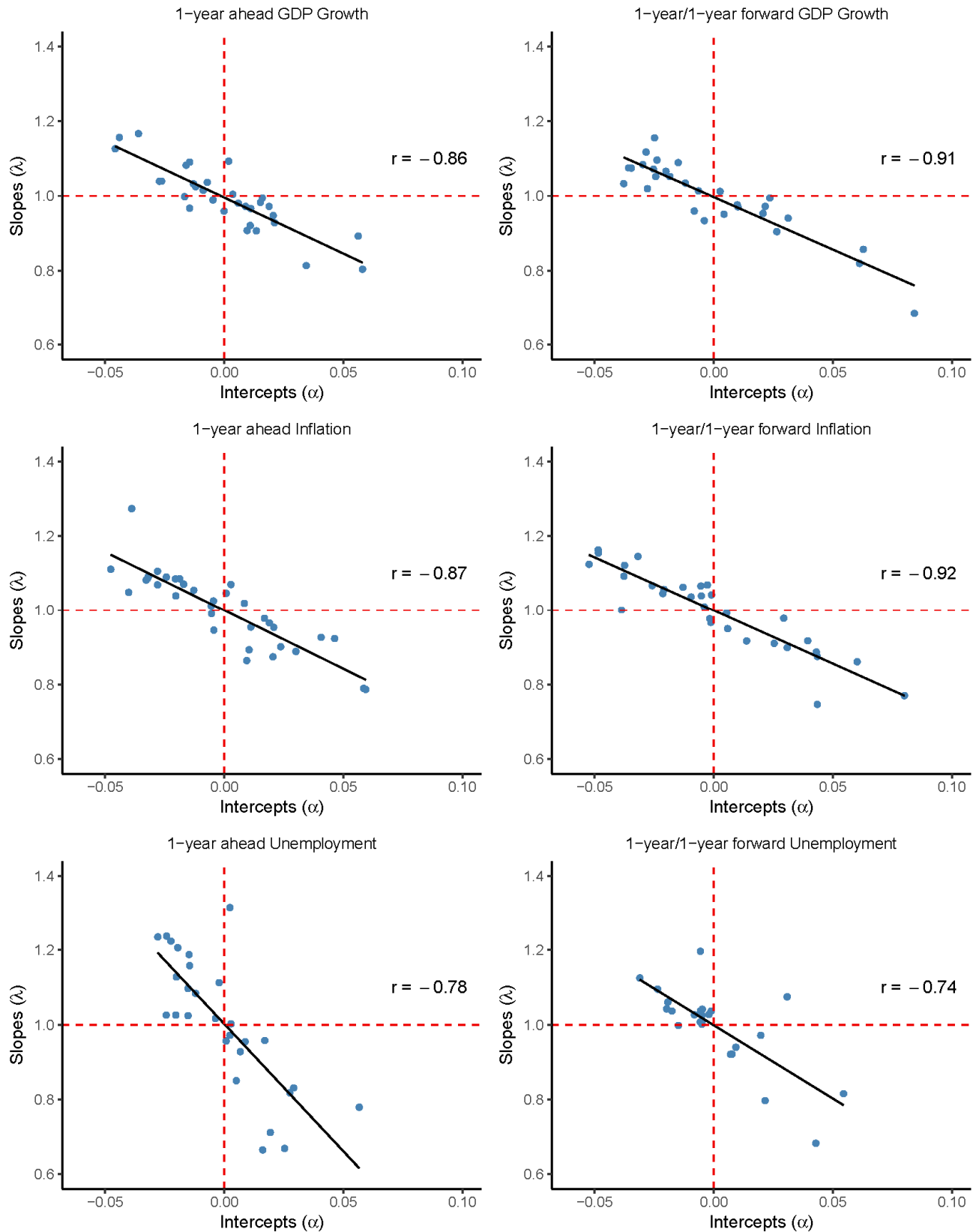


Figure 6

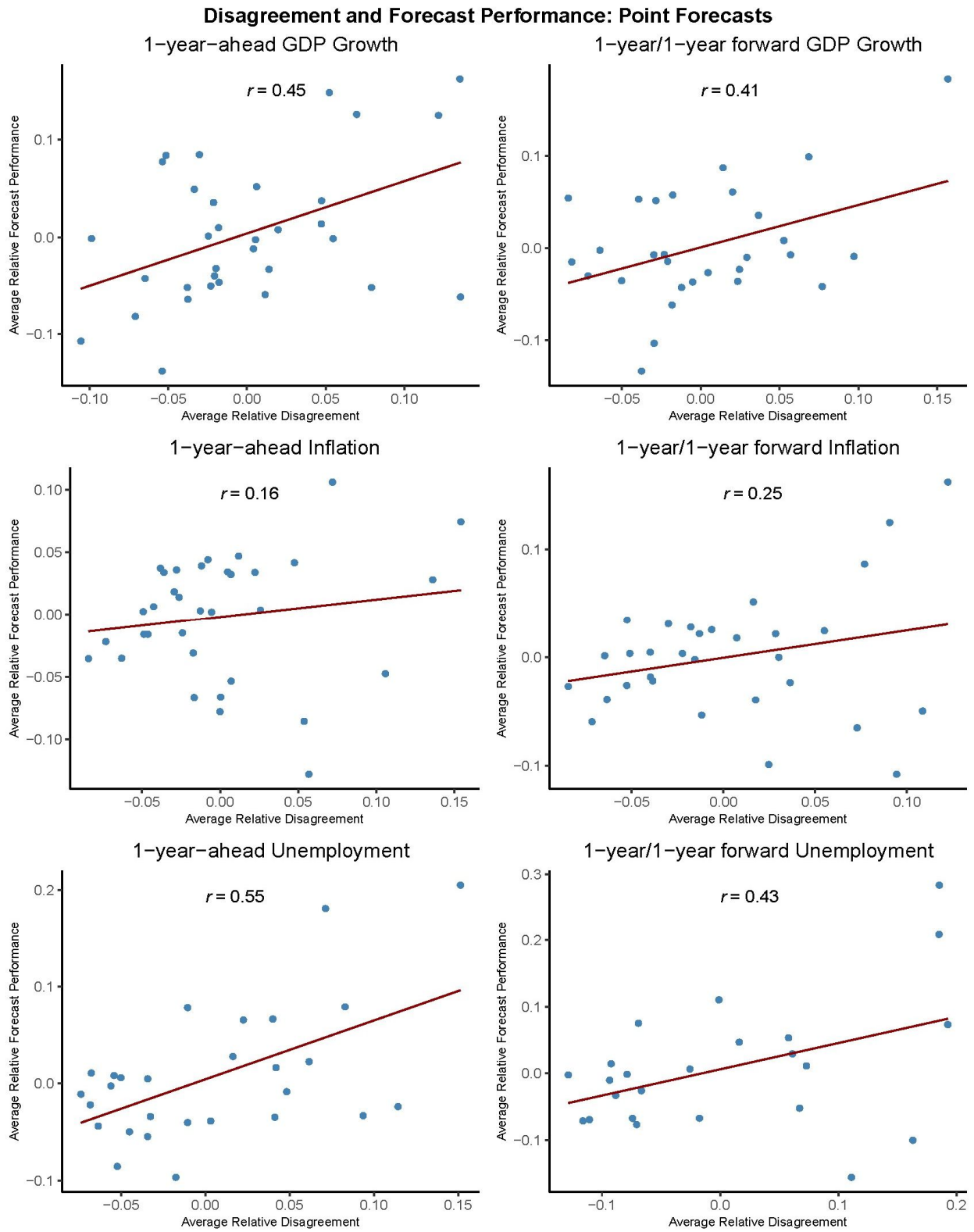


Figure 7

Forecast Performance and Fitted Regression Lines 1-year ahead GDP Growth

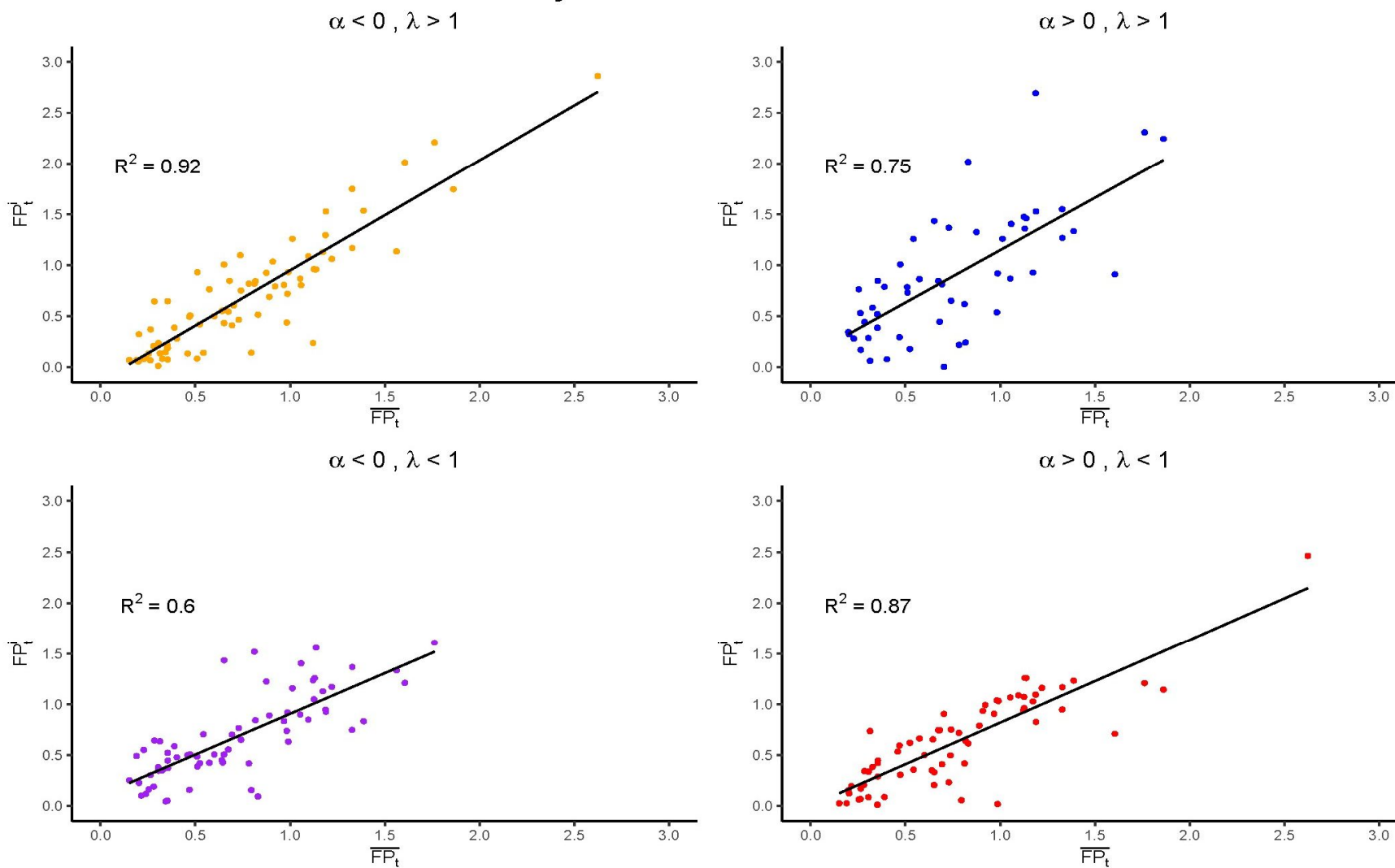
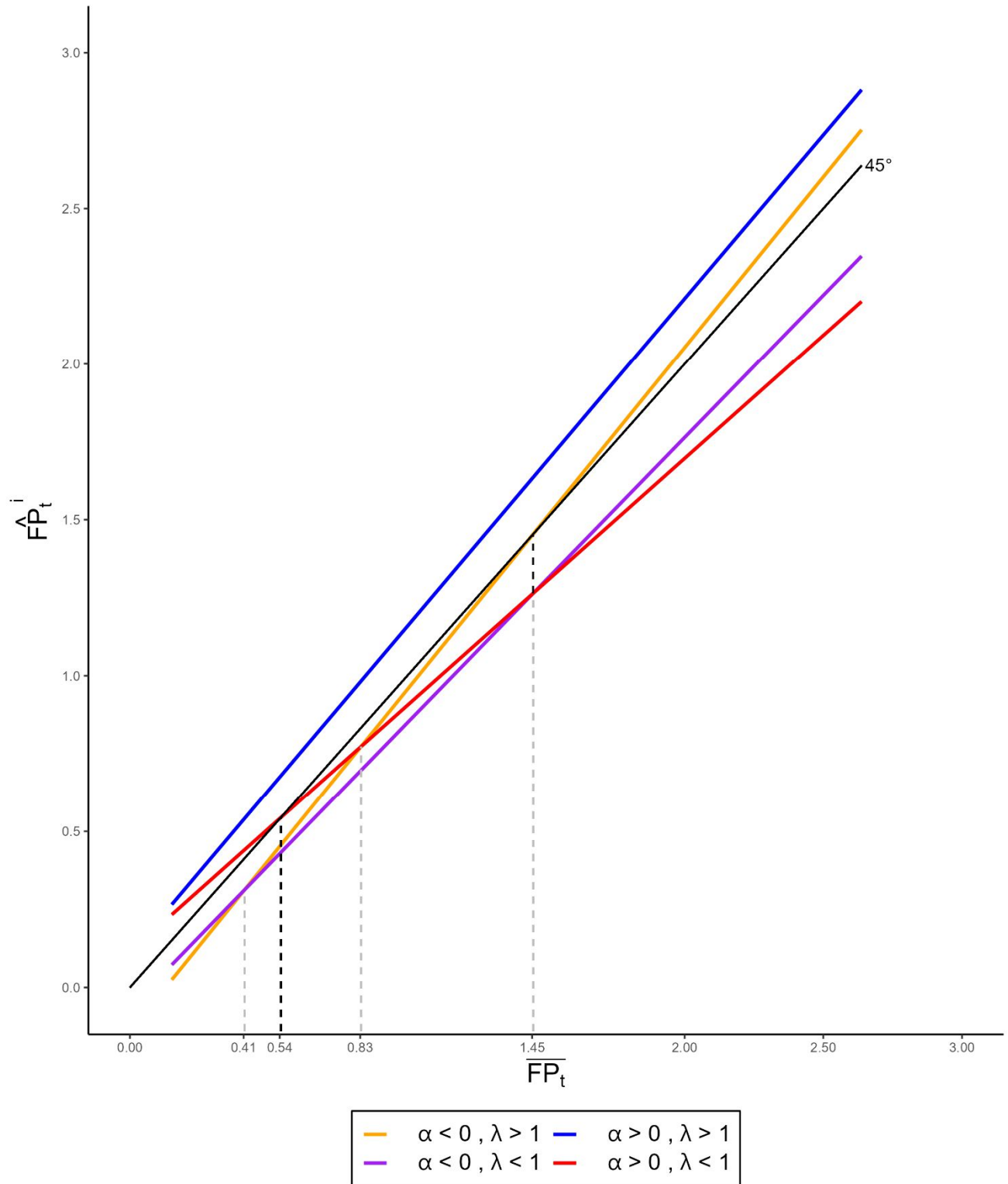


Figure 8

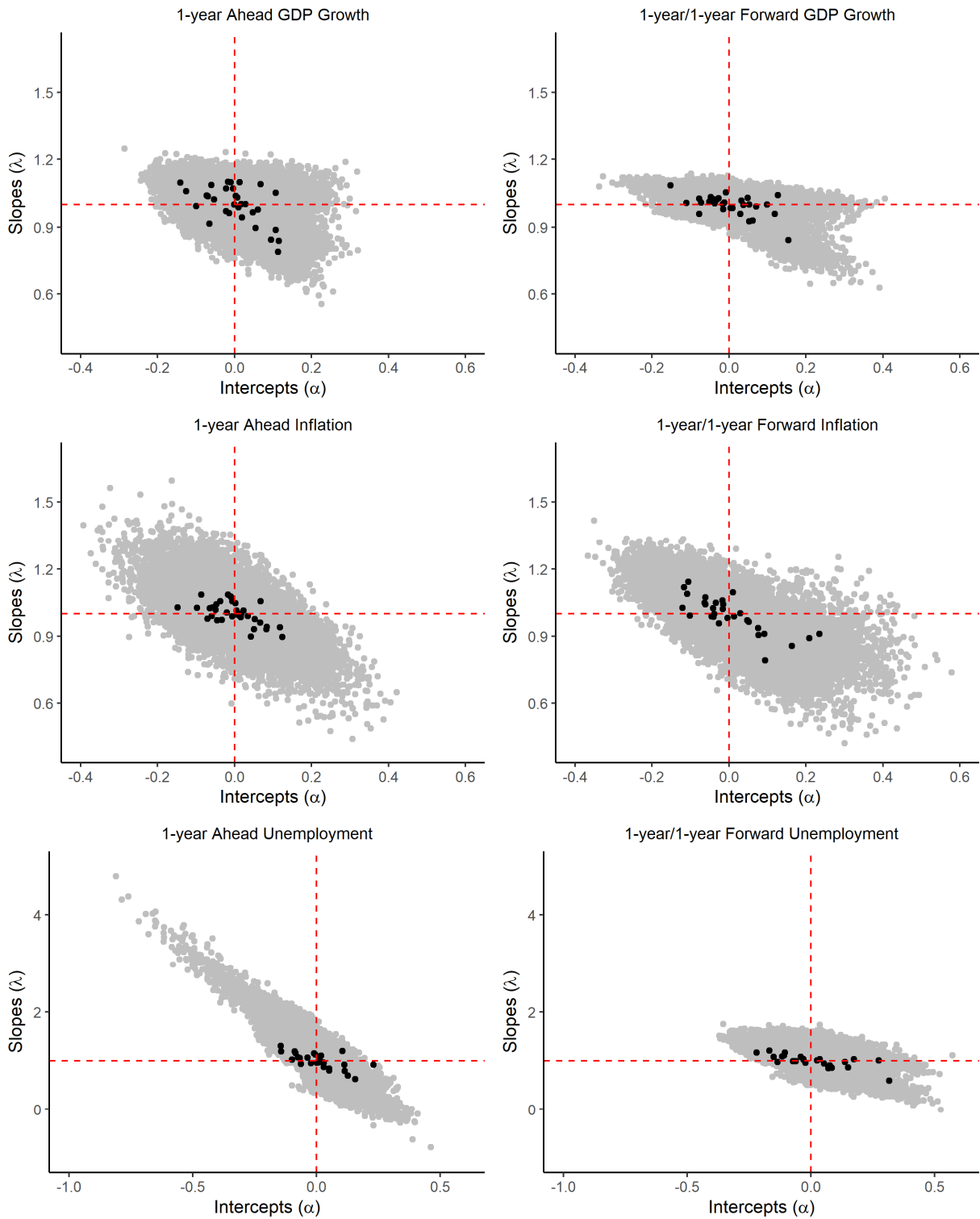
Estimated Forecast Performance Profiles with Crossings
1-year ahead GDP Growth



Grey dashed lines depict crossings of individual forecast performance profiles.
Black dashed lines depict crossings of individual forecast performance profiles with consensus forecast performance.

Figure 9

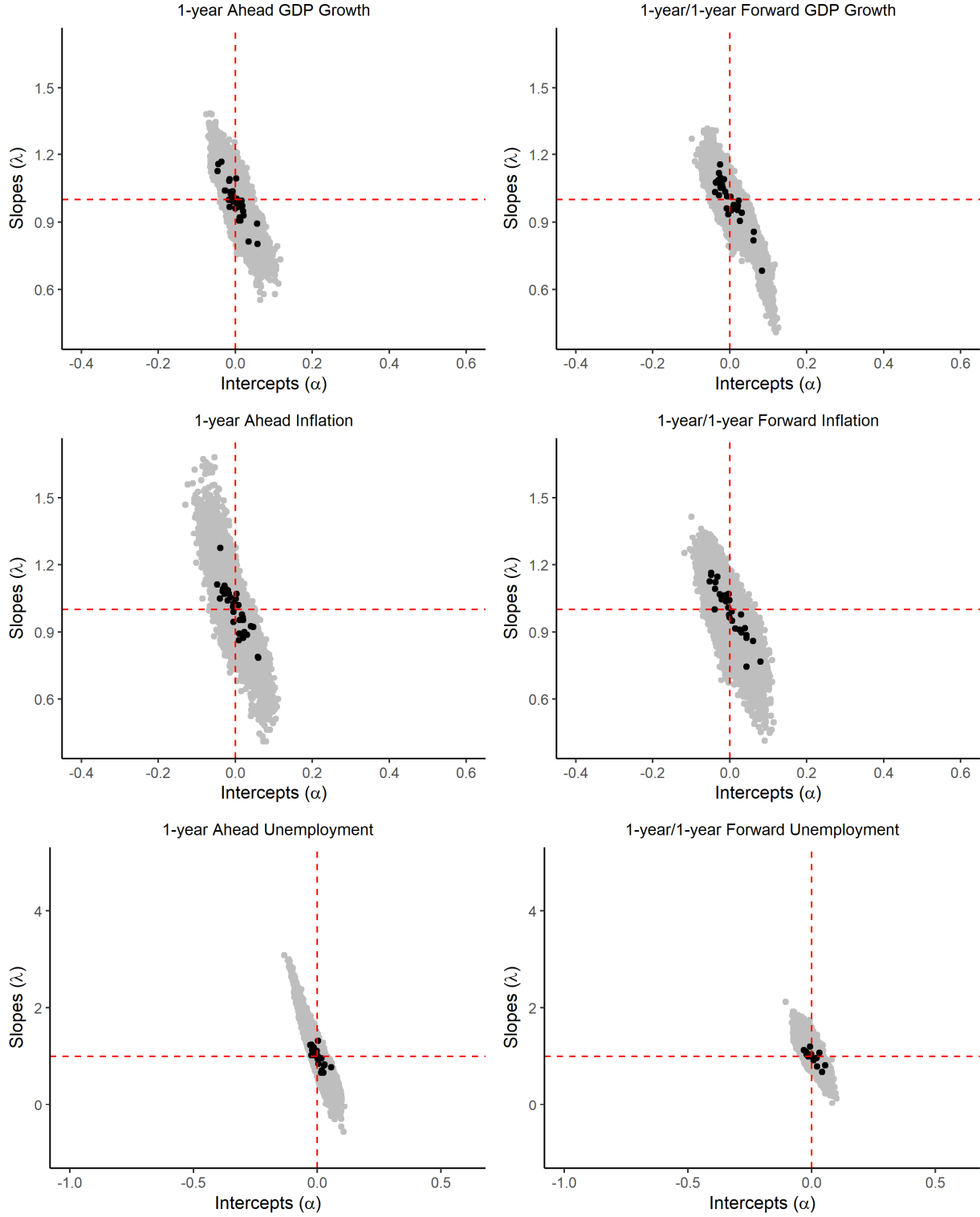
Estimated and Simulated Parameter Pairings: Point Forecasts



Note: Black dots are estimated values and grey dots are simulated values from the estimated joint distributions.

Figure 10

Estimated and Simulated Parameter Pairings: Density Forecasts



Note: Black dots are estimated values and grey dots are simulated values from the estimated joint distributions.

Figure 11

Highest Aggregate Percentage in a Quadrant: Point Forecasts

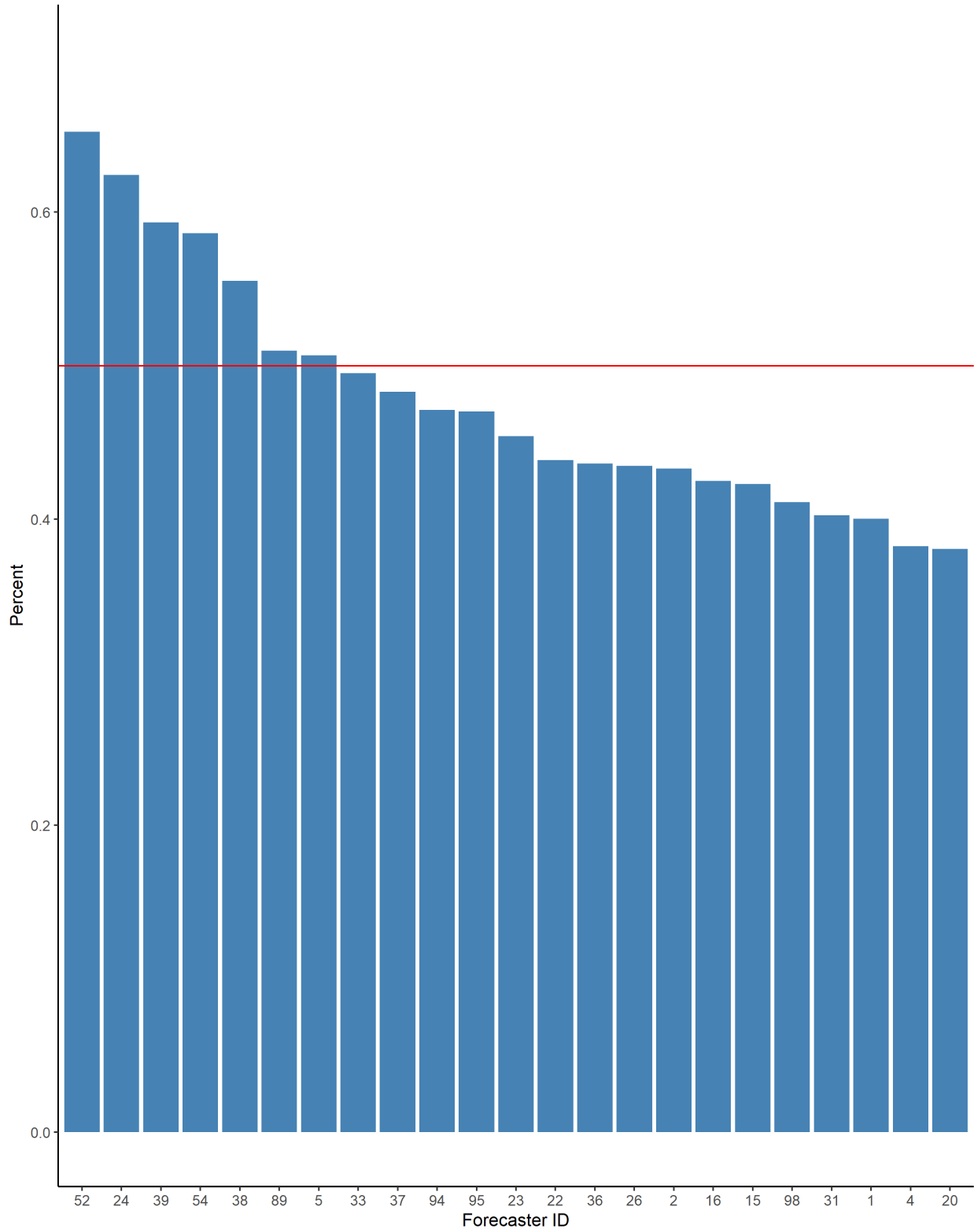


Figure 12

Highest Aggregate Percentage in a Quadrant: Density Forecasts

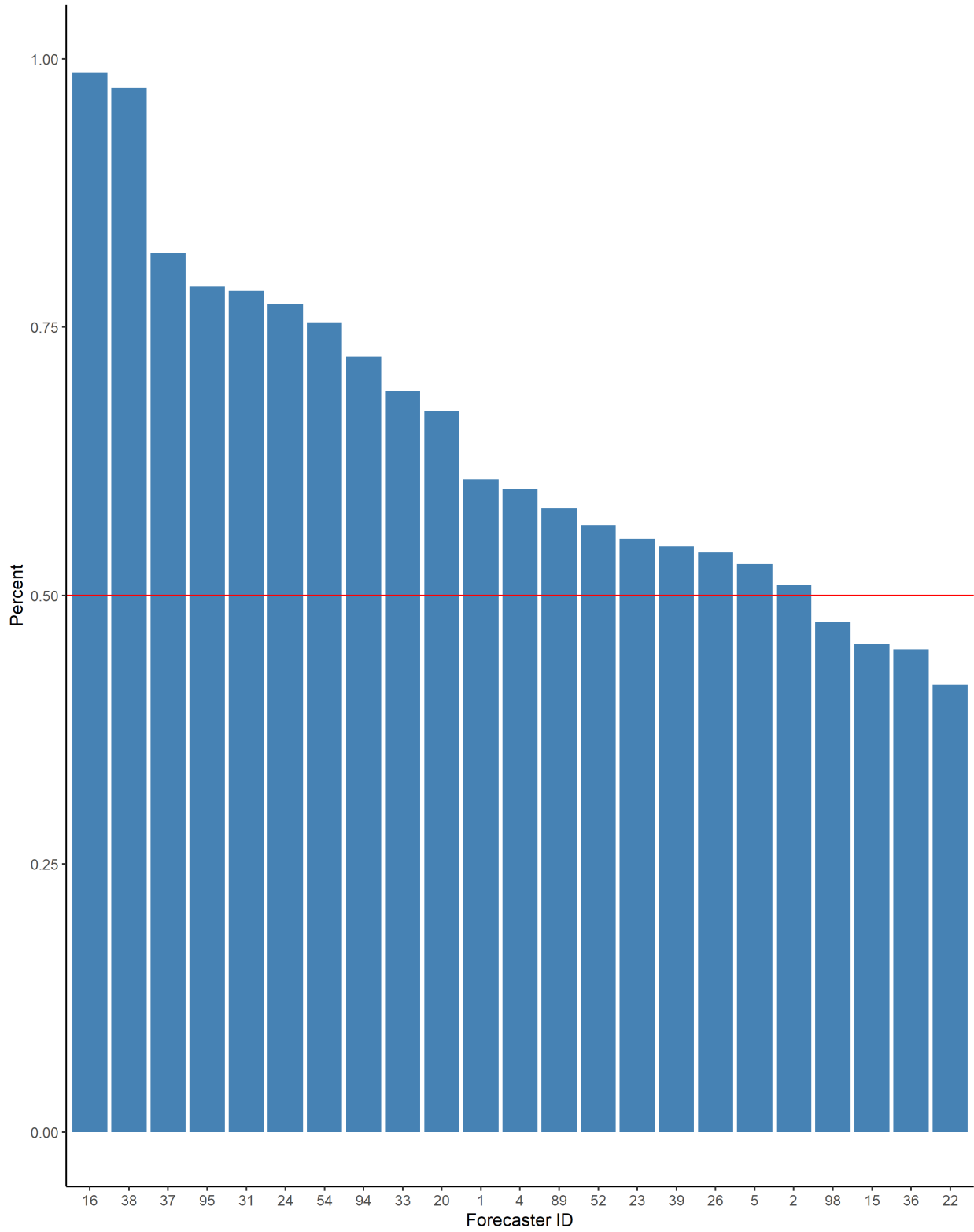
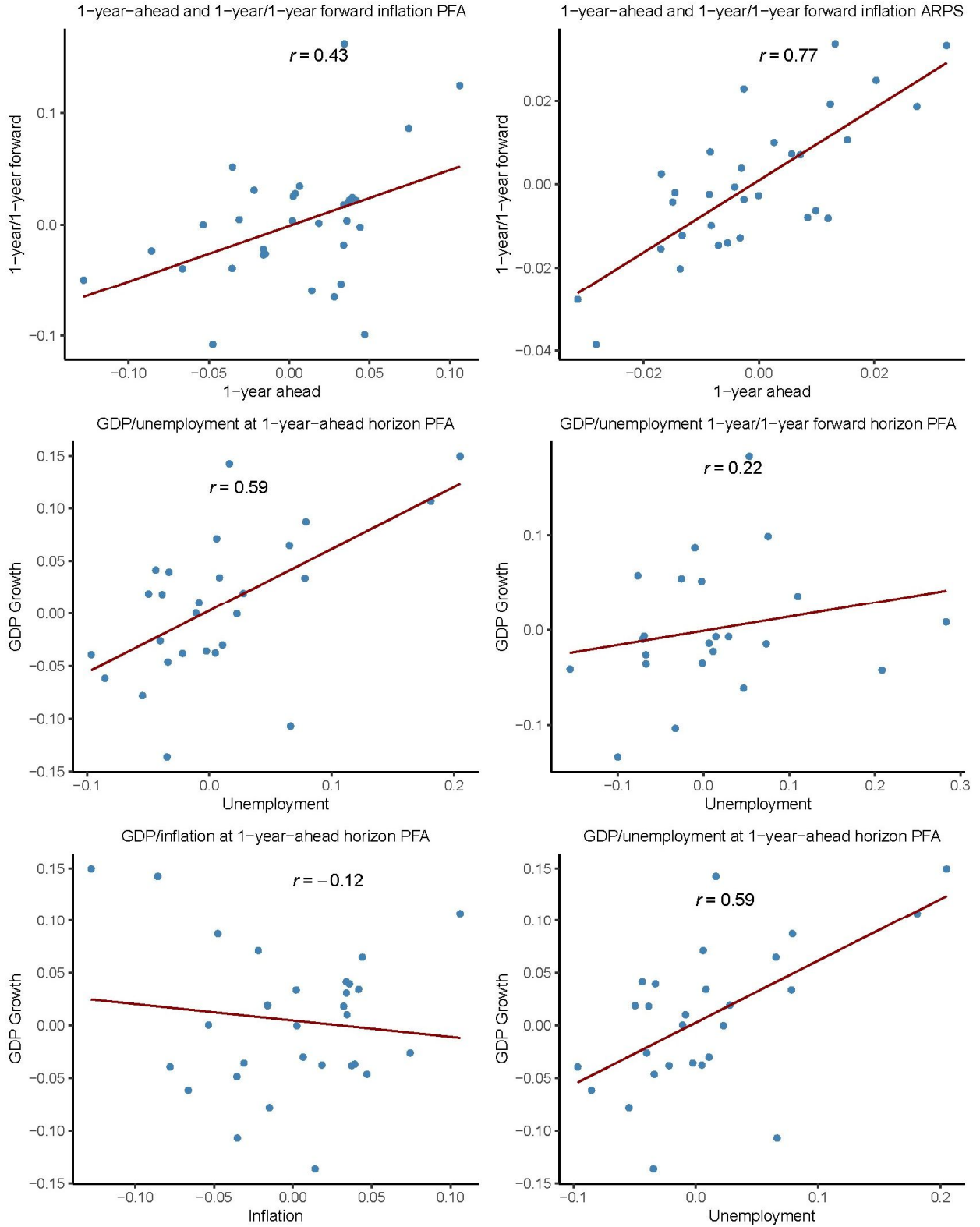


Figure 13

Forecast Performance Comparisons: Data Type, Target Variables, and Horizons



References

- Boero, Gianna, Jeremy Smith, and Kenneth F. Wallis. "The Measurement and Characteristics of Professional Forecasters' Uncertainty." *Journal of Applied Econometrics* 30 (December 2015): 1029-1046.
- Bowles, Carlos, Roberta Friz, Veronique Genre, Geoff Kenny, Aidan Meyler, and Tuomas Rautanen. "The ECB Survey of Professional Forecasters (SPF): A Review After Eight Years' Experience." ECB Occasional Paper No 59. European Central Bank, April 1, 2007.
- Bruine de Bruin, Wandi, Charles F. Manski, Giorgio Topa, and Wilbert van der Klaauw. "Measuring Consumer Uncertainty About Future Inflation." *Journal of Applied Econometrics* 26 (May 2011): 454-478.
- Clements, Michael P. "Forecaster Efficiency and Disagreement: Evidence Using Individual-Level Survey Data." *Journal of Money, Credit and Banking* 54 (April 2022): 537-568.
- Coibion, Olivier, and Yuriy Gorodnichenko. "What Can Survey Forecasts Tell Us About Information Rigidities?" *Journal of Political Economy* 120 (February 2012): 116-159.
- , and Yuriy Gorodnichenko. "Information Rigidity and the Expectations Formation Process: A Simple Framework and New Facts." *American Economic Review* 105 (August 2015): 2644-2678.
- D'Agostino, Antonello, Kieran McQuinn, and Karl Whelan. "Are Some Forecasters Really Better Than Others?" *Journal of Money, Credit, and Banking* 44 (June 2012): 715-732.
- Garcia, Juan A. "An Introduction to the ECB's Survey of Professional Forecasters." ECB Occasional Paper No 8. European Central Bank, 2003.
- Genre, Veronique, Geoff Kenny, Aidan Meyler, and Allan Timmermann. "Combining Expert Forecasts: Can Anything Beat the Simple Average?" *International Journal of Forecasting* 29 (March 2013): 108-121.
- Glas, Alexander, and Matthias Hartman. "Uncertainty Measures from Partially Rounded Probabilistic Forecast Surveys." *Quantitative Economics* 13 (July 2022): 979-1022.
- Hounyo, Ulrich, and Kajal Lahiri. "Are Some Forecasters Really Better Than Others? A Note." *Journal of Money, Credit and Banking* 55 (April 2023): 577-593.
- Kenny, Geoff, Thomas Kostka, and Federico Masera. "How Informative are the Subjective Density Forecasts of Macroeconomists?" *Journal of Forecasting* 33 (April 2014): 163-185.
- , Thomas Kostka, and Federico Masera. "Can Macroeconomists Forecast Risk? Event-Based Evidence from the Euro-Area SPF." *International Journal of Central Banking* 11 (December 2015): 1-46.
- , Thomas Kostka, and Federico Masera. "Density Characteristics and Density Forecast Performance: A Panel Analysis." *Empirical Economics* 48 (May 2015): 1203-1231.
- MacKowiak, Bartosz, and Mirko Wiederholt. "Optimal Sticky Prices Under Rational Inattention." *American Economic Review* 99 (June 2009): 769-803.
- Mankiw, Gregory N., and Ricardo Reis. "Sticky Information Versus Sticky Prices: A Proposal to Replace the New Keynesian Phillips Curve." *Quarterly Journal of Economics* 117 (November 2002): 1295-1328.
- , Ricardo Reis, and Justin Wolfers. "Disagreement about Inflation Expectations." *NBER Macroeconomics Annual* 18 (2003): 209-248.
- Meyler, Aidan. "Forecast Performance in the ECB SPF: Ability or Chance?" Working Paper Series. European Central Bank, 2020.
- Newey, Whitney K., and Kenneth D. West. "A Simple, Positive Semi-Definite, Heteroskedasticity and Autocorrelation Consistent Covariance Matrix." *Econometrica* 55 (May 1987): 703-708.
- Pesaran, M. Hashem. "Estimation and Inference in Large Heterogeneous Panels with a Multifactor Error Structure." *Econometrica* 74 (July 2006): 967-1012.

- Qu, Ritong, Allan Timmermann, and Yinchu Zhu. "Do Any Economists Have Superior Forecasting Skills?" Working Paper. University of California - San Diego, October 1, 2019.
- , Allan Timmermann, and Yinchu Zhu. "Comparing Forecasting Performance in Cross-Sections." *Journal of Econometrics* 237 (December 2021).
- Rich, Robert W., and Joseph Tracy. "A Closer Look at the Behavior of Uncertainty and Disagreement: Micro Evidence from the Euro Area." *Journal of Money, Credit, and Banking* 53 (February 2021): 233-253.
- Sims, Christopher A. "Implications of Rational Inattention." *Journal of Monetary Economics* 50 (April 2003): 665-690.
- Timmerman, Allan. "Forecast Comparisons." In *Handbook of Economic Forecasting, Volume 1 Chapter 4*, edited by G. Elliott, C. Granger and A. Timmerman, 135-196. Elsevier, 2006.
- Woodford, Michael. "Imperfect Common Knowledge and the Effects of Monetary Policy." In *Knowledge, Information, and Expectations in Modern Macroeconomics: In Honor of Edmund S. Phelps*, edited by Philippe Aghion, Roman Frydman, Joseph Stiglitz and Michael Woodford. Princeton, Princeton University Press, 2003.

Appendix

1: Test for Distributional Homogeneity of Predictive Performance Metrics

To demonstrate that the restriction $(\alpha_j, \lambda_j) = (0, 1)$ provides a test for distributional homogeneity of the first and second moments of forecast performance metrics across survey participants, we begin by taking expectations of equation (3) in the paper:

$$\begin{aligned} E[FP_{t+h|t}^j] &= E[\alpha_j + \lambda_j (\overline{FP_{t+h|t}}) + \varepsilon_{t+h|t}^j] \\ E[FP_{t+h|t}^j] &= \alpha_j + \lambda_j E(\overline{FP_{t+h|t}}) \end{aligned} \tag{1.1}$$

As discussed in Section II, we define interchangeability as the condition that there are no systematic differences across participants in their forecast behavior. Equating (1.1) for participant i and j , we have:

$$E(FP_{t+h|t}^i) = E(FP_{t+h|t}^k) \tag{1.2}$$

Substituting (1.1) into (1.2) yields:

$$\alpha_i + \lambda_i E(\overline{FP_{t+h|t}}) = \alpha_k + \lambda_k E(\overline{FP_{t+h|t}}) \tag{1.3}$$

where the absence of systematic differences across forecasters in their predictive performance requires that $\alpha_i = \alpha_k = \alpha$ and $\lambda_i = \lambda_k = \lambda$.

Using (1.1) and the condition that $\alpha_i = \alpha_k = \alpha$ and $\lambda_i = \lambda_k = \lambda$ from (1.3) allows us to derive the following expression for the expectation of average forecast performance:

$$\begin{aligned} E[\overline{FP_{t+h|t}}] &= E\left(\frac{1}{N_t} \left[\sum_{j=1}^{N_t} FP_{t+h|t}^j \right]\right) \\ &= \frac{1}{N_t} E\left[\sum_{j=1}^{N_t} \left(\alpha_j + \lambda_j E(\overline{FP_{t+h|t}}) + \varepsilon_{t+h|t}^j\right)\right] \\ &= \frac{1}{N_t} \left[\sum_{j=1}^{N_t} \left(\alpha_j + \lambda_j E(\overline{FP_{t+h|t}})\right) \right] \\ &= \frac{1}{N_t} \left[\sum_{N_t} \left(\alpha + \lambda E(\overline{FP_{t+h|t}})\right) \right] \\ &= \alpha + \lambda E(\overline{FP_{t+h|t}}) \end{aligned} \tag{1.4}$$

where N_t denotes the number of survey respondents at time t and where (1.4) only holds if $\alpha = 0$ and $\lambda = 1$.

Looking further at equation (1.1) and abstracting from the discussion of α_j , the condition that $\lambda_j = 1$ requires covariance $\left(FP_{t+h|t}^j, \overline{FP_{t+h|t}} \right)$ to be equal to the variance $\left(\overline{FP_{t+h|t}} \right)$. This would be similar for all other respondents. Because the ε 's in equation (1.1) are assumed to be uncorrelated across respondents because the common correlated effects (CCE) estimator controls for cross-sectional dependence, then covariance $\left(FP_{t+h|t}^j, \overline{FP_{t+h|t}} \right)$ will be equal to variance $\left(\varepsilon_{t+h|t}^j \right)$, which will only equal variance $\left(\overline{FP_{t+h|t}} \right)$ if all the variances across respondents are equal. Consequently, the restriction $(\alpha_j, \lambda_j) = (0, 1)$ provides a test that:

$$E(FP_{t+h|t}^1) = E(FP_{t+h|t}^2) = \dots = E(FP_{t+h|t}^N) \cap V(FP_{t+h|t}^1) = V(FP_{t+h|t}^2) = \dots = V(FP_{t+h|t}^N) \quad (1.5)$$

which involves an equality of the mean and variance of the distribution of the forecast performance metrics across all respondents.

2: Hounyo – Lahiri (2023) Testing Procedure

The following tables report the results from applying the Hounyo-Lahiri (2023) testing procedure for equal predictive performance to point forecasts and density forecasts from the European Central Bank Survey of Professional Forecasters. Following Hounyo and Lahiri (2023), the testing procedure for the point forecasts was implemented by using the normalized squared error statistic given by:

$$\overline{POINT} FP_{t+h|t}^j = \frac{(X_{t+h} - E_t^j[X_{t+h}])^2}{\frac{1}{N_t} \sum_{j=1}^{N_t} (X_{t+h} - E_t^j[X_{t+h}])^2} \quad (2.1)$$

where N_t is again the number of survey respondents at time t and the remaining notation is defined in equation (6) of the paper.

For the density forecasts, we implemented the testing procedure using the normalized absolute rank probability score given by:

$$\overline{DENSITY} FP_{t+h|t}^j = \frac{DENSITY FP_{t+h|t}^j}{\frac{1}{N_t} \sum_{j=1}^{N_t} DENSITY FP_{t+h|t}^j} \quad (2.2)$$

where the non-negative domain of the absolute rank probability score obviates the application of an additional operator on the forecast performance metric and the remaining notation is defined in equation (7) of the paper.

As shown, the results in Table 1A (point forecasts) and Table 2A (density forecasts) strongly reject the null hypothesis of comparable forecast performance and document that the performance of forecasters across the various percentiles is more accurate than what can be explained by random chance (generated from their wild bootstrap procedure) at conventional levels of significance. The overall results and p -values are very similar to those reported by Hounyo and Lahiri (2023).

Following Hounyo and Lahiri (2023) we also implemented the tests using a restrictive data set, where we excluded forecasters who scored worse than the 80th percentile. Table 3A (point forecasts) and Table 4A (density forecasts) report the results from this exercise. As shown, excluding the forecasters does not change the conclusions.

Table 1A – Point Forecast Data

		Best	5	25	50	75	Worst
Forecaster performance		0.669	0.685	0.905	0.995	1.115	1.547
GDP growth: one-year-ahead	$(Q_5^*; Q_{95}^*)$	(1.11; 1.70)	(1.22; 1.83)	(1.58; 2.33)	(1.76; 2.59)	(1.97; 2.91)	(2.71; 4.21)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.705	0.719	0.897	0.973	1.095	1.474
GDP growth: one-year/one-year forward	$(Q_5^*; Q_{95}^*)$	(1.20; 1.79)	(1.27; 1.88)	(1.59; 2.32)	(1.76; 2.59)	(2.00; 2.92)	(2.55; 4.05)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.690	0.696	0.862	0.987	1.126	1.485
HICP inflation: one-year-ahead	$(Q_5^*; Q_{95}^*)$	(1.13; 1.69)	(1.24; 1.81)	(1.54; 2.25)	(1.79; 2.59)	(2.02; 2.94)	(2.67; 4.04)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.754	0.784	0.834	0.940	1.135	1.686
HICP inflation: one-year/one-year forward	$(Q_5^*; Q_{95}^*)$	(1.40; 2.10)	(1.45; 2.18)	(1.62; 2.39)	(1.86; 2.71)	(2.20; 3.20)	(3.33; 5.32)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.626	0.658	0.847	0.954	1.104	2.091
Unemployment: one-year-ahead	$(Q_5^*; Q_{95}^*)$	(1.18; 1.75)	(1.25; 1.83)	(1.67; 2.45)	(1.89; 2.77)	(2.17; 3.20)	(4.02; 6.18)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.739	0.774	0.844	0.945	1.086	1.908
Unemployment: one-year/one-year forward	$(Q_5^*; Q_{95}^*)$	(1.33; 2.04)	(1.39; 2.13)	(1.65; 2.45)	(1.85; 2.72)	(2.15; 3.20)	(3.63; 5.77)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000

NOTE: This table provides the empirical distribution of forecaster performance for point forecasts of GDP growth, HICP inflation, and unemployment from the European Central Bank Survey of Professional Forecasters. We measure forecast performance for point forecasts using the average of the normalized squared forecast error in equation (1.1). The figures in parentheses $(Q_5^*; Q_{95}^*)$ refer to the 5th and the 95th percentiles generated by the cross-sectional and serial correlation bootstrap procedure of Hounyo and Lahiri (2023). The reported p -value is the proportion of the 999 bootstrap replications that are less than observed forecast performance.

Table 2A – Density Forecast Data

		Best	5	25	50	75	Worst
Forecaster performance		0.806	0.812	0.942	0.964	1.079	1.309
GDP growth: one-year-ahead	$(Q_5^*; Q_{95}^*)$	(1.38; 2.07)	(1.44; 2.13)	(1.63; 2.39)	(1.76; 2.59)	(1.92; 2.81)	(2.32; 3.44)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.806	0.834	0.921	0.969	1.076	1.260
GDP growth: one-year/one-year forward	$(Q_5^*; Q_{95}^*)$	(1.49; 2.26)	(1.60; 2.37)	(1.81; 2.62)	(1.94; 2.80)	(2.11; 3.07)	(2.47; 3.67)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.768	0.820	0.920	0.971	1.084	1.261
HICP inflation: one-year-ahead	$(Q_5^*; Q_{95}^*)$	(1.48; 2.14)	(1.62; 2.39)	(1.81; 2.60)	(1.94; 2.78)	(2.14; 3.06)	(2.49; 3.64)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.761	0.786	0.907	0.991	1.068	1.285
HICP inflation: one-year/one-year forward	$(Q_5^*; Q_{95}^*)$	(1.70; 2.53)	(1.78; 2.64)	(2.07; 2.99)	(2.23; 3.24)	(2.45; 3.53)	(2.87; 4.23)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.727	0.765	0.924	0.969	1.057	1.471
Unemployment: one-year-ahead	$(Q_5^*; Q_{95}^*)$	(1.42; 2.08)	(1.49; 2.21)	(1.79; 2.60)	(1.91; 2.79)	(2.09; 3.05)	(2.86; 4.22)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.866	0.870	0.930	0.981	1.015	1.323
Unemployment: one-year/one-year forward	$(Q_5^*; Q_{95}^*)$	(1.71; 2.61)	(1.77; 2.64)	(1.93; 2.87)	(2.08; 3.05)	(2.19; 3.24)	(2.79; 4.20)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000

NOTE: This table provides the empirical distribution of forecaster performance for density forecasts of GDP growth, HICP inflation, and unemployment from the European Central Bank Survey of Professional Forecasters. We measure forecast performance for density forecasts using the average of the normalized absolute rank probability score in equation (1.2). The figures in parentheses $(Q_5^*; Q_{95}^*)$ refer to the 5th and the 95th percentiles generated by the cross-sectional and serial correlation bootstrap procedure of Hounyo and Lahiri (2023). The reported p -value is the proportion of the 999 bootstrap replications that are less than observed forecast performance.

Table 3A – Point Forecast Data: Restricted to Best 80%

		Best	5	25	50	75	Worst
Forecaster performance		0.696	0.707	0.957	1.038	1.103	1.207
GDP growth: one-year-ahead	$(Q_5^*; Q_{95}^*)$	(1.18; 1.78)	(1.29; 1.89)	(1.65; 2.39)	(1.87; 2.66)	(2.01; 2.89)	(2.31; 3.53)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.763	0.773	0.913	1.007	1.097	1.215
GDP growth: one-year/one-year forward	$(Q_5^*; Q_{95}^*)$	(1.14; 1.75)	(1.20; 1.80)	(1.45; 2.16)	(1.60; 2.38)	(1.75; 2.59)	(2.02; 3.21)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.739	0.748	0.915	1.002	1.107	1.275
HICP inflation: one-year-ahead	$(Q_5^*; Q_{95}^*)$	(1.21; 1.84)	(1.31; 1.93)	(1.61; 2.35)	(1.81; 2.62)	(2.01; 2.93)	(2.42; 3.63)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.834	0.840	0.890	0.972	1.120	1.279
HICP inflation: one-year/one-year forward	$(Q_5^*; Q_{95}^*)$	(1.52; 2.30)	(1.57; 2.36)	(1.72; 2.55)	(1.93; 2.71)	(2.18; 3.23)	(2.70; 4.17)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.695	0.711	0.932	1.038	1.067	1.268
Unemployment: one-year-ahead	$(Q_5^*; Q_{95}^*)$	(1.32; 1.95)	(1.37; 2.02)	(1.79; 2.64)	(2.01; 2.92)	(2.20; 3.20)	(2.63; 4.08)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.788	0.805	0.920	0.962	1.137	1.236
Unemployment: one-year/one-year forward	$(Q_5^*; Q_{95}^*)$	(1.42; 2.17)	(1.48; 2.23)	(1.75; 2.59)	(1.92; 2.82)	(2.19; 3.32)	(2.54; 3.84)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000

NOTE: This table provides the empirical distribution of forecaster performance for point forecasts of GDP growth, HICP inflation, and unemployment from the European Central Bank Survey of Professional Forecasters. We measure forecast performance for point forecasts using the average of the normalized squared forecast error in equation (1.1). The figures in parentheses $(Q_5^*; Q_{95}^*)$ refer to the 5th and the 95th percentiles generated by the cross-sectional and serial correlation bootstrap procedure of Hounyo and Lahiri (2023). The reported p -value is the proportion of the 999 bootstrap replications that are less than observed forecast performance.

Table 4A – Density Forecast Data: Restricted to Best 80%

		Best	5	25	50	75	Worst
Forecaster performance		0.832	0.839	0.950	0.992	1.055	1.150
GDP growth: one-year-ahead	$(Q_5^*; Q_{95}^*)$	(1.30; 1.90)	(1.36; 1.94)	(1.52; 2.17)	(1.65; 2.34)	(1.76; 2.51)	(1.93; 2.79)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.838	0.858	0.953	0.988	1.055	1.152
GDP growth: one-year/one-year forward	$(Q_5^*; Q_{95}^*)$	(1.39; 2.16)	(1.48; 2.23)	(1.66; 2.48)	(1.76; 2.62)	(1.88; 2.81)	(2.05; 3.08)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.804	0.840	0.949	0.989	1.053	1.151
HICP inflation: one-year-ahead	$(Q_5^*; Q_{95}^*)$	(1.54; 2.25)	(1.63; 2.44)	(1.85; 2.71)	(1.97; 2.86)	(2.10; 3.05)	(2.36; 3.45)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.791	0.812	0.934	1.021	1.072	1.141
HICP inflation: one-year/one-year forward	$(Q_5^*; Q_{95}^*)$	(1.64; 2.48)	(1.70; 2.55)	(1.96; 2.91)	(2.12; 3.13)	(2.24; 3.33)	(2.46; 3.70)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.780	0.806	0.973	1.017	1.052	1.158
Unemployment: one-year-ahead	$(Q_5^*; Q_{95}^*)$	(1.51; 2.21)	(1.59; 2.31)	(1.87; 2.73)	(1.98; 2.89)	(2.09; 3.05)	(2.31; 3.47)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000
Forecaster performance		0.892	0.901	0.942	1.019	1.040	1.075
Unemployment: one-year/one-year forward	$(Q_5^*; Q_{95}^*)$	(1.67; 2.53)	(1.72; 2.57)	(1.88; 2.77)	(2.01; 2.95)	(2.11; 3.10)	(2.24; 3.32)
	p-value	0.000	0.000	0.000	0.000	0.000	0.000

NOTE: This table provides the empirical distribution of forecaster performance for density forecasts of GDP growth, HICP inflation, and unemployment from the European Central Bank Survey of Professional Forecasters. We measure forecast performance for density forecasts using the average of the normalized absolute rank probability score in equation (1.2). The figures in parentheses $(Q_5^*; Q_{95}^*)$ refer to the 5th and the 95th percentiles generated by the cross-sectional and serial correlation bootstrap procedure of Hounyo and Lahiri (2023). The reported p -value is the proportion of the 999 bootstrap replications that are less than observed forecast performance.

3: Figures Displaying the Full Data Set for 1-year-ahead GDP Growth Forecasts

Figure 7A and Figure 8A exclude an outlier observation for three of the four participants associated with realized GDP growth in 2008:Q3 or 2008:Q4. As discussed in the main text, this exclusion was done for presentational purposes. The following figures include the full range of \overline{FP} values. The R^2 values reported in Figure 7A are identical to those that appear in Figure 7 in the paper.

Figure 7A

Forecast Performance and Fitted Regression Lines
1-year ahead GDP Growth

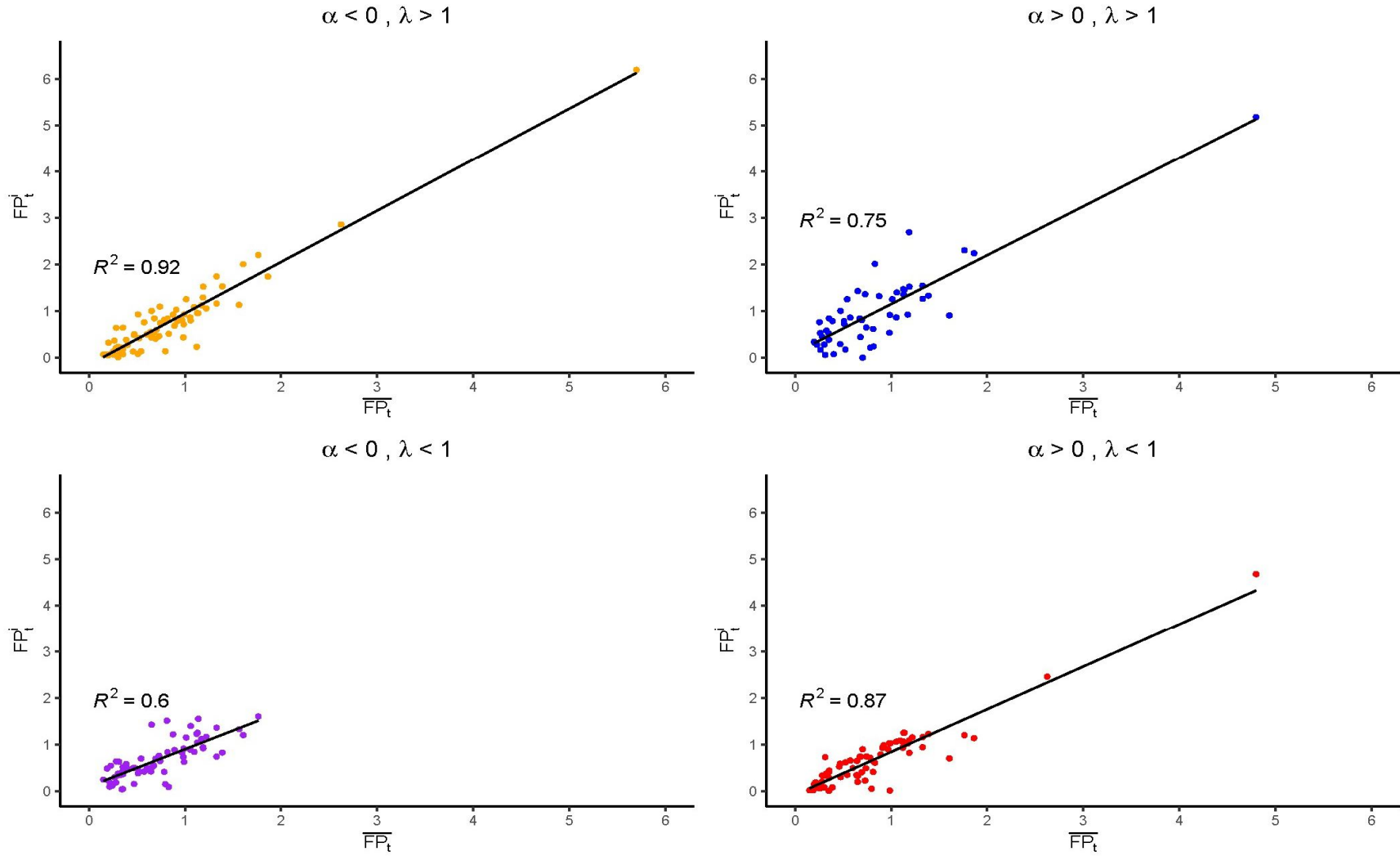
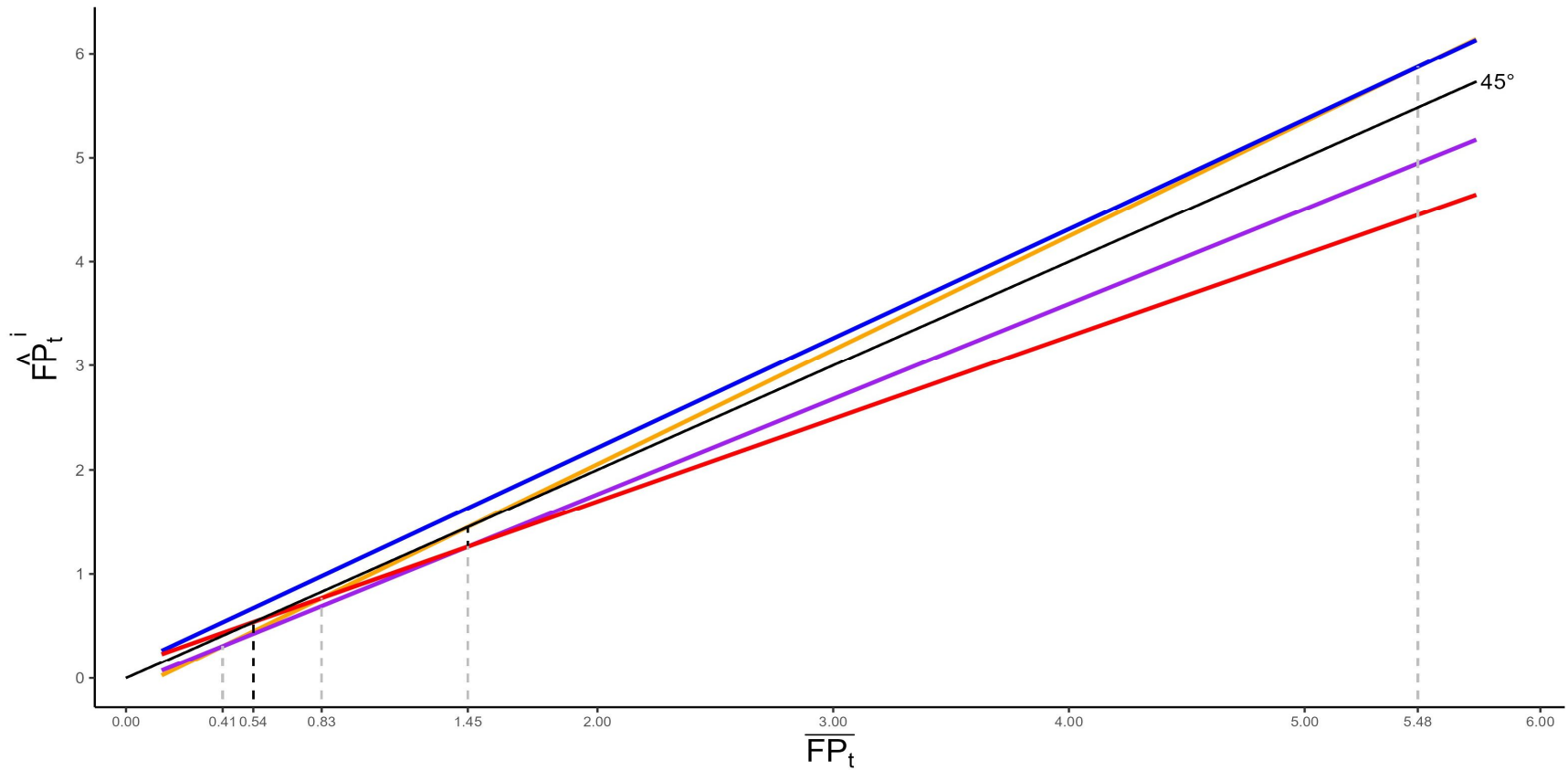


Figure 8A

Estimated Forecast Performance Profiles with Crossings 1-year ahead GDP Growth



— $\alpha < 0, \lambda > 1$	— $\alpha > 0, \lambda > 1$
— $\alpha < 0, \lambda < 1$	— $\alpha > 0, \lambda < 1$

Grey dashed lines depict crossings of individual forecast performance profiles.
Black dashed lines depict crossings of individual forecast performance profiles with consensus forecast performance.

4: Forecast Combination Exercise

For this exercise, we assume knowledge of the difficulty of the forecasting environment (\overline{FP}) as well as the estimated parameters and the root mean square error (RMSE) of equation (3) at the individual level. For each forecaster j , we calculate the probability of the forecaster being more accurate than the cross-sectional average as:

$$\hat{p}_{t+h|t}^j = \Phi \left(\frac{\overline{FP}_{t+h|t} - \overline{FP}_{t+h|t}^j}{\hat{\sigma}^j} \right) \quad (4.1)$$

where Φ is the cumulative standard normal distribution and $\hat{\sigma}^j$ is the RMSE for forecaster j .

Using Figure 8 as a reference, it is relatively straightforward to understand the intuition for equation (4.1). The numerator on the right-hand side of equation (4.1) is the vertical distance between the 45-degree line and a respondent's predicted performance for a particular survey date, which is then normalized by the respondent's RMSE. If the respondent's predicted performance coincides with average predictive performance, the associated probability is 50 percent. Values of \overline{FP} that lie above (below) \overline{FP} are associated with individual predictive performance that is relatively less (more) accurate than average performance and results in probability values less (more) than 50 percent, with further deviations between \overline{FP} and \overline{FP} generating larger changes in probability.

The performance weighted consensus forecast is given by:

$$FP_{t+h|t}^{Performance-weighted} = \sum_{j=1}^{N_t} w_{t+h|t}^j FP_{t+h|t}^j \quad (4.2)$$

where the weight for forecaster j , $w_{t+h|t}^j$, is:

$$w_{t+h|t}^j = \frac{\hat{p}_{t+h|t}^j}{\sum_{i=1}^{N_t} \hat{p}_{t+h|t}^i} \quad (4.3)$$

We carry out the exercise for the one-year-ahead GDP growth forecasts using eight surveys that are closest to the \overline{FP} values in Table 3. As shown in Table 5A, we find that the performance-weighted consensus forecast outperforms the (equally weighted) consensus forecast in each of these eight surveys. Moreover, the median performance improvement is 24 percent, with the smallest improvement at 12.6 percent and the largest improvement at 35.5 percent.

Table 5A

<p>Consensus Forecast</p> \overline{FP}	<p>Performance-weighted Forecast</p> $\sum_{i=1}^{N_t} w_{t+h t}^i \overline{FP}_{t+h t}^i$	<p>Ratio</p> $\frac{\sum_{i=1}^{N_t} w_{t+h t}^i \overline{FP}_{t+h t}^i}{\overline{FP}}$
0.2536	0.2056	0.811
0.5094	0.3742	0.735
0.7408	0.6328	0.854
0.9907	0.7530	0.760
1.562	1.1760	0.753
1.861	1.2005	0.645
4.798	4.1952	0.874
5.698	4.2086	0.738