

Comparing Predictive Accuracy in the Presence of a Loss Function Shape Parameter*

Sander Barendse Andrew J. Patton
Oxford University Duke University

December 8, 2020

Abstract

We develop tests for out-of-sample forecast comparisons based on loss functions that contain shape parameters. Examples include comparisons using average utility across a range of values for the level of risk aversion, comparisons of forecast accuracy using characteristics of a portfolio return across a range of values for the portfolio weight vector, and comparisons using recently-proposed “Murphy diagrams” for classes of consistent scoring rules. An extensive Monte Carlo study verifies that our tests have good size and power properties in realistic sample sizes, particularly when compared with existing methods which break down when the number of values considered for the shape parameter grows. We present three empirical illustrations of the new test.

Keywords: Forecasting, model selection, out-of-sample testing, nuisance parameters.

J.E.L. codes: C53, C52, C12.

*We thank Dick van Dijk, Erik Kole, Chen Zhou, and seminar participants at Oxford University for valuable discussions and feedback. The first author also acknowledges financial support from the Erasmus Trustfonds. All errors remain our own. Correspondence to: Sander Barendse, Nuffield College, New Road, Oxford, UK. Email address: sander.barendse@economics.ox.ac.uk.

1 Introduction

Forecast comparison problems in economics and finance invariably rely on a loss function or utility function, and in many cases these functions contain a shape parameter, for example, when comparing two forecasting models based on average utility. In many cases, there is no single specific value for the shape parameter that is of interest, rather there is a *range* of values that are of interest to the researcher. In this case a null hypothesis of superior predictive accuracy has a continuum of testable implications, but most existing work instead tests predictive accuracy at a few ad hoc values of the shape parameter. This paper combines work on forecast comparison tests, see Diebold and Mariano (1995), White (2000), and Giacomini and White (2006), with bootstrap theory for empirical processes, see Bühlmann (1995), to provide forecast comparison tests that allow for inference across a range of values of a loss function shape parameter.

A leading example of an application that depends on a shape parameter is a test of equal expected utility in which the utility function is parameterized by a risk aversion parameter. Given that economists have not converged on what value of risk aversion is appropriate (see, e.g., Bliss and Panigirtzoglou, 2004, for discussion) it is desirable to consider a range of reasonable risk aversion values, instead of testing at some single value. Current practice usually evaluates the hypothesis of equal expected utility at one or a select few risk aversion parameter values, see Fleming et al. (2001), Marquering and Verbeek (2004), Engle and Colacito (2006) and DeMiguel et al. (2007) for example.

Another application that involves a continuum of testable implications is when one evaluates forecasts from multivariate models on the basis of their implied forecasts of univariate quantities, for example, evaluating a multivariate volatility model through its forecasts of Value-at-Risk for portfolios of the underlying assets. In practice, comparisons of quantile forecasts of portfolio returns, as generated from multivariate models, often only consider the equal weighted portfolio or some other fixed combination of portfolio constituents, see McAleer and Da Veiga (2008), Santos et al. (2012) and Kole et al. (2017) for example. Considering only a single weight vector can fail to reveal the sensitivity, or the robustness, of the ranking of two models to the choice of weight vector.

A final, recent example is the comparison of forecasts using “elementary scoring rules,” see Ehm et al. (2016). These authors show that the family of loss functions (or “scoring rules”) that are consistent for a given statistical functional (e.g., the mean, a quantile, an expectile) can be

represented as a convex combination of elementary scoring rules. The plot of the elementary scoring rules is called a “Murphy diagram,” and Ehm et al. (2016) note that joint testing of the Murphy diagram is not yet fully developed. Ziegel et al. (2020) introduce tests for Murphy diagrams based on controlling the family-wise error rate but such corrections can perform poorly in large-scale multiple testing problems, see Hand (1998) and White (2000). Moreover, their tests consider only a finite subset of the testable implications instead of a continuum.

This paper develops new out-of-sample tests for multiple testing problems over a continuum of shape parameters, including the above three examples, which do not rely on bounds such as the Bonferroni correction, and which take into account the time series nature of the data used in most forecasting settings. To the best of our knowledge, such tests have not been considered in the literature to date. We consider tests of equal predictive ability (two-sided tests) and superior predictive ability (one-sided tests), and provide a framework for forecast dominance between two models based on one-sided tests that use opposing null hypotheses. We derive our tests using functions of the Diebold-Mariano test statistic for each value of the shape parameter in its range, and obtain critical values using the moving blocks bootstrap of Bühlmann (1995), which is applicable to weakly dependent empirical processes indexed by classes of functions, and is general enough to cover our cases of interest: loss functions parameterized by a vector that can take values in a bounded subset of Euclidian space.

Our tests build on the out-of-sample testing framework of Diebold and Mariano (1995) and West (1996). Similar to Giacomini and White (2006) we consider evaluating the forecasting *method*, which, in addition to the forecasting model, includes the estimation scheme and choices of in-sample and out-of-sample periods. Our problem has some similarities with the work of White (2000), Hansen (2005) and Hubrich and West (2010) who present predictive ability tests to compare a benchmark model with finitely many alternative models using a single loss function. In contrast, we consider only two models but we compare them using a continuum of loss functions. As such, our testing problem cannot be addressed using the work of those papers.

Recently, Jin et al. (2017) and Post et al. (2018) introduced multiple comparison tests for general loss functions over relatively large classes of loss functions (although the loss function should take its minimum value when the forecast error is zero), by translating the multiple hypothesis into hypotheses about stochastic dominance and “nonparametric forecast optimality,” respectively. This allows the authors to develop tests based on the empirical distribution func-

tion or empirical likelihood. Our paper differs from this strand of research in two ways. First, we consider classes of loss functions that are economically motivated, whereas the previous authors abstract away from this choice. In many settings, such as the ones considered in this paper, it is realistic to assume the researcher has at least some knowledge of what constitutes a relevant loss, and thus what values of the shape parameter are relevant. Ruling out economically uninteresting loss functions from the general set of loss functions can improve test power. Second, our paper applies to scenarios in which the loss function is not minimal at zero forecast error, or where the very notion of a forecast error is unclear. Such loss functions include, for instance, many utility functions as well as the Fissler and Ziegel (2016) loss function for Value-at-Risk and Expected Shortfall.

We show via an extensive simulation study that the proposed testing methods have good size and power properties in realistic sample sizes. These positive results stand in stark contrast to the two most familiar existing methods used to compare forecasts across a range of loss function parameter values: we find that the Wald test has finite-sample size as high as 50% for a 5% level test, while tests based on a Bonferroni correction are conservative and suffer, as a result, from low power.

We consider three empirical applications of the proposed new tests. Firstly, in a comparison of expected utility of equal weighted and minimum-variance portfolio strategies (see, e.g., DeMiguel et al., 2007) we show that our tests are able to reject the null hypothesis when existing alternative methods cannot. Secondly, we consider a tests of portfolio quantile forecasts generated by multivariate GARCH-DCC and RiskMetrics models and we find that our tests again are able to detect violations of the null hypothesis where existing methods cannot. Finally, we consider tests based on the Murphy diagrams for quantile forecasts generated by GARCH and RiskMetrics models, an application where existing methods simply cannot be applied due to the nature of the testing problem.

The remainder of the paper is structured as follows. In Section 2 we discuss our three illustrative examples. In Section 3 we present the general testing framework and develop our tests. In Section 4 we use Monte Carlo experiments to study the small sample properties of our tests in settings close to our illustrative examples. In Section 5 we explore these settings empirically. Section 6 concludes.

2 Loss function shape parameters in practice

We consider three representative examples of forecast comparison scenarios. In all of these examples we consider loss differences defined as

$$L_{t+1}(\gamma) = S(Y_{t+1}, g_t^{\mathcal{A}}(\gamma); \gamma) - S(Y_{t+1}, g_t^{\mathcal{B}}(\gamma); \gamma) \quad (1)$$

where S is the loss function (or scoring rule), $\gamma \in \Gamma \subset \mathbb{R}^d$ is the shape parameter of the loss function, Y_{t+1} is the target variable, and $g_t^i(\gamma)$ is the forecast of Y_{t+1} made using model $i \in \{\mathcal{A}, \mathcal{B}\}$, which may depend on the shape parameter γ .

We focus on tests of uniform superior predictive ability, which consider the hypotheses:

$$H_0 : E[L_{t+1}(\gamma)] \leq 0 \quad \forall \gamma \in \Gamma \quad (2)$$

$$\text{vs. } H_1 : E[L_{t+1}(\gamma^\dagger)] \geq \Delta > 0 \text{ for some } \gamma^\dagger \in \Gamma \quad (3)$$

We will assume that lower values of the scoring rule are preferred, and so H_0 implies that model \mathcal{A} is weakly better than model \mathcal{B} for all $\gamma \in \Gamma$, while H_1 implies that model \mathcal{B} is strictly better than model \mathcal{A} for some $\gamma \in \Gamma$. In the supplemental appendix we also consider a test of uniform equal predictive accuracy:

$$H'_0 : E[L_{t+1}(\gamma)] = 0 \quad \forall \gamma \in \Gamma \quad (4)$$

$$\text{vs. } H'_1 : \left| E[L_{t+1}(\gamma^\dagger)] \right| \geq \Delta > 0 \text{ for some } \gamma^\dagger \in \Gamma \quad (5)$$

2.1 Comparisons based on expected utility

Two forecasting models each generate forecasts of optimal portfolio weights and we seek to compare them in terms of out-of-sample average utility from the resulting portfolio returns. The portfolio returns are obtained as $Y'_{t+1} g_t^i(\gamma)$, where Y_{t+1} is a vector of returns on the underlying assets, and $g_t^i(\gamma)$ is the forecasted optimal portfolio weights from model $i \in \{\mathcal{A}, \mathcal{B}\}$ assuming a preference parameter γ , and per-period utility is computed using some utility function $u(\cdot; \gamma)$. For instance, when $u(\cdot; \gamma)$ is the exponential utility function γ denotes the (scalar) risk aversion

parameter, and $\Gamma = [a, b]$, for some $0 < a < b < \infty$. In this case we set

$$S(Y_{t+1}, g_t^i(\gamma); \gamma) = -u(Y_{t+1}' g_t^i(\gamma); \gamma), \text{ for } i \in \{\mathcal{A}, \mathcal{B}\} \quad (6)$$

so that lower values of the scoring rule indicate better performance.

2.2 Multivariate forecast comparison based on portfolio characteristics

Let Y_{t+1} denote some vector of returns, and compute portfolio returns as $\tilde{Y}_{t+1}(\gamma) = \gamma' Y_{t+1}$ for some weight vector γ . We are interested in forecasting some statistic ψ_t of $\tilde{Y}_{t+1}(\gamma)$ for all portfolio weights $\gamma \in \Gamma$. When we consider all long-only portfolios with weights summing to one, the parameter space Γ is the unit simplex. If we consider α -quantile forecasts of the portfolio return, then we may use the “tick” loss function to measure forecast performance and so set

$$S(Y_{t+1}, g_t^i(\gamma, \alpha); \gamma, \alpha) = (\mathbf{1}\{\gamma' Y_{t+1} \leq g_t^i(\gamma, \alpha)\} - \alpha)(g_t^i(\gamma, \alpha) - \gamma' Y_{t+1}), \text{ for } i \in \{\mathcal{A}, \mathcal{B}\}, \quad (7)$$

where $\mathbf{1}\{\cdot\}$ equals one if the argument is true and zero otherwise.

2.3 Forecast comparisons via Murphy diagrams

Let Y_{t+1} denote some scalar return, and let ξ_t denote some statistic of Y_{t+1} , such as a mean or quantile. If ξ_t is elicitable, see Gneiting (2011a), then there exists a family of scoring rules (loss functions), \mathcal{S} such that for any scoring rule $S^* \in \mathcal{S}$ it holds

$$E[S^*(Y_{t+1}, \xi_t)] \leq E[S^*(Y_{t+1}, x)], \quad \forall x \in \mathcal{X} \quad (8)$$

where \mathcal{X} is the support of ξ_t . That is, ξ_t has lower expected loss than any other forecast, x , for all scoring rules $S^* \in \mathcal{S}$. The scoring rule S^* is then said to be “consistent” for the statistic ξ_t . Many statistics, such as the mean, quantile, and expectile, admit *families* of consistent scoring functions, see Gneiting (2011a). For example, the mean is well-known to be elicitable using the quadratic loss function, and more generally it is elicitable using any “Bregman” loss function. The α -quantile is elicitable using the tick loss function in equation (7), and more generally using any “generalized piecewise linear” (GPL) loss function, see Gneiting (2011b). Comparisons of forecasts are usually done using a *single* consistent scoring rule (e.g., using mean squared error to compare estimates of the mean), however Patton (2020) shows that in

the presence of parameter estimation error or misspecified models the ranking of forecasts can be sensitive to the specific scoring rule used.

In a recent paper, Ehm et al. (2016) show that any scoring rule that is consistent for a quantile or an expectile (the latter nesting the mean as a special case) can be represented as:

$$S^*(Y_{t+1}, x) = \int_{-\infty}^{\infty} \tilde{S}(Y_{t+1}, x; \gamma) dH(\gamma) \quad (9)$$

for some non-negative measure H , and $\gamma \in \Gamma \subset \mathbb{R}$, where \tilde{S} is an “elementary” scoring rule, defined below. A plot of the average elementary scores across all values of γ is called a “Murphy diagram” by Ehm et al. (2016). If one forecast’s Murphy diagram lies below that of another forecast for all values of γ , then it has lower average loss for *any* consistent scoring rule.

With the above representation of a consistent scoring rule as a mixture of elementary scoring rules, we can consider rankings across *all* consistent scoring rules, overcoming the sensitivity discussed in Patton (2020). For example, to compare α -quantile forecasts we set

$$S(Y_{t+1}, g_t^i(\alpha); \gamma, \alpha) = (\mathbf{1}\{Y_{t+1} < g_t^i\} - \alpha)(\mathbf{1}\{\gamma < g_t^i(\alpha)\} - \mathbf{1}\{\gamma < Y_{t+1}\}), \text{ for } i \in \{\mathcal{A}, \mathcal{B}\} \quad (10)$$

and then test for forecast superiority across all $\gamma \in \mathbb{R}$. The right-hand side of the above equation is the quantile elementary scoring rule from Ehm et al. (2016).

3 Forecast comparison tests in the presence of a loss function shape parameter

Consider the stochastic process $W = \{W_t : \Omega \rightarrow \mathbb{R}^{N+s}, N \in \mathbb{N}_+, s \in \mathbb{N}, t = 1, 2, \dots\}$ defined on a complete probability space (Ω, \mathcal{F}, P) . We partition the observed vector W_t as $W_t = (Y_t, X_t)$, where $Y_t : \Omega \rightarrow \mathbb{R}^N$ is a vector a variables of interest and $X_t : \Omega \rightarrow \mathbb{R}^s$ is a vector of explanatory variables. We define $\mathcal{F}_t = \sigma(W_1, \dots, W_t)$.

To fix notation, we let $|A| = (\text{tr}(A'A))^{1/2}$ denote the Euclidean norm of a matrix A , and $\|A\|_q = (E|A|^q)^{1/q}$ denote the \mathcal{L}_q norm of a random matrix. Finally, \Rightarrow denotes weak convergence with respect to the uniform metric.

We denote the total sample size by T and the out-of-sample size by n . We consider moving or fixed window forecasts generated with in-sample periods of size m , such that the forecast for period $t + 1$ is obtained using observations at periods $t - m + 1, \dots, t$ with the moving scheme,

and $1, \dots, m$ with the fixed scheme, respectively.

We consider some (scalar) measurable loss difference:

$$L_{t+1}(\gamma) = L(W_{t+1}, W_t, \dots, W_{t-m+1}; \gamma) \quad (11)$$

that takes as arguments $m + 1 < \infty$ elements of W and some parameter vector $\gamma \in \Gamma \subset \mathbb{R}^d$ that is independent of W , with Γ a bounded set. As in Giacomini and White (2006), our assumption that $m < \infty$ imposes a limited memory condition on the forecasting methods, which precludes methods with model parameters estimated over expanding windows, but allows for those estimated over fixed and rolling windows of finite length. Note that McCracken (2020) cautions that in some applications additional assumptions are needed to apply the theory in Giacomini and White (2006) when a fixed-window estimation scheme is used. (We use a rolling window estimation scheme in our empirical analysis.) Also note that our Giacomini-White style analysis of forecast performance rules out inference on forecast performance at the pseudo-true parameter value, as in West (1996) and Clark and McCracken (2001).

3.1 Superior predictive ability tests

We will focus attention on tests of uniform superior predictive ability, which consider the null and alternative hypotheses, H_0 and H_1 , defined in equations (2)-(3). The case of equal predictive accuracy, which considers the hypotheses H'_0 and H'_1 defined in equations (4)-(5), is perhaps less economically interesting than that of superior accuracy, and so we relegate further discussion of that case to the supplemental appendix. Notice that H'_0 in equation (4) is the element of H_0 least favorable to the alternative, and is the point at which we derive the limit distribution of our test statistic. Also note that these null hypotheses involve unconditional expected loss, rather than expected loss conditional on some other variable, which is also considered in Giacomini and White (2006).

To develop a test of H_0 we consider the Diebold and Mariano (1995) test statistic as a function of $\gamma \in \Gamma$ and then take the supremum of that function over Γ :

$$t_n(\gamma) \equiv \sqrt{n} \frac{\bar{L}_n(\gamma)}{\hat{\sigma}_n(\gamma)} \quad (12)$$

$$\sup t_n \equiv \sup_{\gamma \in \Gamma} t_n(\gamma). \quad (13)$$

where $\bar{L}_n(\gamma) \equiv \frac{1}{n} \sum_{t=m}^{T-1} L_{t+1}(\gamma)$, and $\hat{\sigma}_n^2(\gamma)$ denotes a consistent estimator of $\sigma^2(\gamma) \equiv E[L_{t+1}(\gamma)^2]$.

It should be noted that when autocorrelation is present in $L_{t+1}(\gamma)$, $t_n(\gamma)$ does not converge in distribution to a standard normal limit, because $\hat{\sigma}_n^2(\gamma)$ is not a heteroskedasticity and autocorrelation consistent (HAC) estimator of the asymptotic covariance matrix of $\sqrt{n}\bar{L}_n(\gamma)$, see, e.g. Newey and West (1987)). We are not aware of strong uniform law of large numbers results for HAC estimators, which are required in our theoretical results below. As noted by Hansen (2005), $\hat{\sigma}_n^2(\gamma)$ is not required to be a consistent estimator of the variance of $\sqrt{n}\bar{L}_n(\gamma)$, because the bootstrap accounts for time series features in the data to obtain critical values of our tests. Indeed, in some scenarios it might be better not to studentize at all, and fix $\hat{\sigma}_n^2(\gamma) = 1$ instead. Such scenarios include those for which $\hat{\sigma}_n^2(\gamma)$ may be close to zero in small samples.

The test statistic $\sup t_n$ can be written as a function $v(t_n)$, where v maps functionals on Γ to \mathbb{R} and we write $t_n = \{t_n(\gamma) : \gamma \in \Gamma\}$ as a random function on Γ . The function v is continuous with respect to the uniform metric, monotonic in the sense that if $Z_1(\gamma) \leq Z_2(\gamma)$ for all γ then $v(Z_1) \leq v(Z_2)$, and has the property that if $Z(\gamma) \rightarrow \infty$ for some γ then $v(Z) \rightarrow \infty$. Adopting this notation also facilitates easily handling of test statistics for tests of equal predictive ability, which are discussed in the supplemental appendix.

We derive the asymptotic distribution of the test statistic above using the following assumptions.

Assumption 1. $\{W_t\}$ is stationary and β -mixing (absolutely regular), with $\beta(t) = c_\beta a^t$, for some finite constant c_β , and $0 < a < 1$.

Assumption 2. $E[\sup_{\gamma \in \Gamma} |L_{t+1}(\gamma)|^{4r}] < \infty$, for some $r > 1$, and for all t .

Assumption 3. $\|L_{t+1}(\gamma) - L_{t+1}(\gamma')\|_{4r} \leq C|\gamma - \gamma'|^\lambda$, for some $C < \infty$, $\lambda > 0$, and for all $\gamma, \gamma' \in \Gamma$, and t .

Assumption 4. $\hat{\sigma}_n^2(\gamma) \xrightarrow{a.s.} \sigma^2(\gamma)$ uniformly over $\gamma \in \Gamma$. Moreover, $\inf_{\gamma \in \Gamma} \sigma^2(\gamma) > 0$.

The β -mixing condition in Assumption 1, which is stronger than α -mixing, but weaker than ϕ -mixing, is usually assumed when deriving functional CLTs for time series data with unbounded absolute moments. Bühlmann (1995) notes that the mixing rate is satisfied for ARMA(p, q) processes with innovations dominated by the Lebesgue measure. Boussama et al. (2011) provides conditions under which multivariate GARCH models satisfy geometric ergodicity, and Bradley et al. (2005, Thm. 3.7) shows that geometric ergodicity implies β -mixing with

at least exponential rate, satisfying Assumption 1. Assumption 2 is a standard moment condition, and Assumption 3 requires a Lipschitz condition to hold for these moments. Assumption 4 requires that $\hat{\sigma}_n^2(\cdot)$ satisfies a strong Uniform Law of Large Numbers (see, e.g., Andrews, 1992), and also imposes uniform non-singularity of $\sigma^2(\cdot)$.

As a building block, we first provide a functional CLT for a demeaned version of $t_n(\gamma)$: $\tau_n(\gamma) \equiv \sqrt{n}(\bar{L}_n(\gamma) - E[L_1(\gamma)])/\hat{\sigma}_n(\gamma)$, as H_0 allows for nonpositive $E[L_1(\cdot)]$. We then provide inference results for the test statistic $v(t_n)$.

Theorem 1. *Let Assumptions 1 to 4 be satisfied. It follows that $\sqrt{n}(\bar{L}_n(\cdot) - E[L_1(\cdot)]) \Rightarrow Z(\cdot)$, for some Gaussian process $Z(\cdot)$ with covariance kernel $\Sigma(\cdot, \cdot) \equiv \lim_{n \rightarrow \infty} Cov(\sqrt{n}\bar{L}_n(\cdot), \sqrt{n}\bar{L}_n(\cdot))$. Moreover, $v(\tau_n) \xrightarrow{d} v(\tilde{t})$, with $\tilde{t}(\cdot) \equiv Z(\cdot)/\sigma(\cdot)$.*

The following result establishes inference under the null and alternative hypotheses.

Theorem 2. *Let Assumptions 1 to 4 be satisfied, and let $\Sigma(\cdot, \cdot)$ be nondegenerate. Under H_0 it follows that $\limsup_{n \rightarrow \infty} P(v(t_n) > c(1 - \alpha)) \leq \lim_{n \rightarrow \infty} P(v(\tau_n) > c(1 - \alpha)) = \alpha$, where $c(1 - \alpha)$ is chosen such that $P(v(\tilde{t}) > c(1 - \alpha)) = \alpha$. Under H_1 it follows that $P(v(t_n) > c(1 - \alpha)) \rightarrow 1$.*

The result in Theorem 2 shows that our test has size of at most α under the one-sided null H_0 . Under the two-sided null H_0' we achieve the nominal rate α . The tests have power approaching one against fixed alternatives. In the case of i.i.d. loss differences $L_{t+1}(\cdot)$ one can employ the results in Andrews and Shi (2013) to obtain nominal size for the one-sided test as well. Unfortunately, these results cannot straightforwardly be extended to loss differences with time-series properties, such as those encountered in the forecasting applications in this paper, and we do not pursue such an extension here.

We establish the consistency of the block bootstrap for general empirical processes of Bühlmann (1995) for $v(t_n)$. The block bootstrap was first studied by Künsch (1989) for general stationary observations. The bootstrap counterpart of $\sqrt{n}\bar{L}_n(\gamma)$ is given by

$$\sqrt{n}\bar{L}_n^*(\gamma) \equiv \sqrt{n} \frac{1}{n} \sum_{t=m}^{T-1} (L_{t+1}^*(\gamma) - \mu_n^*(\gamma)), \quad (14)$$

where $\mu_n^*(\gamma) \equiv \frac{1}{n-l+1} \sum_{i=1}^{n-l+1} \frac{1}{l} \sum_{t=m+i}^{m+i+l-1} L_t(\gamma)$ denotes the expectation of $\frac{1}{n} \sum_{t=m}^{T-1} L_{t+1}^*(\gamma)$ conditional on the original sample, l denotes the block length, and $L_{t+1}^*(\gamma)$ denotes the bootstrap counterpart of $L_{t+1}(\gamma)$. Similarly, let $t_n^*(\gamma) \equiv \sqrt{n}\bar{L}_n^*(\gamma)/\hat{\sigma}_n(\gamma)$, and let $c_n^*(1 - \alpha)$ denote the α -quantile of $t_n^*(\gamma)$.

We impose the following condition on the rate that $l = l(n) \rightarrow \infty$, as $n \rightarrow \infty$.

Assumption 5. *The block length l satisfies $l(n) = O(n^{1/2-\varepsilon})$, for some $0 < \varepsilon < 1/2$.*

The following result establishes consistency of the bootstrap.

Theorem 3. *Let Assumptions 1 to 5 be satisfied. It follows that $\sqrt{n}\bar{L}_n^*(\cdot) \Rightarrow Z(\cdot)$ almost surely.*

Moreover, under H_0 $\limsup_{n \rightarrow \infty} P(v(t_n) > c_n^(1 - \alpha)) \leq \lim_{n \rightarrow \infty} P(v(\tau_n) > c_n^*(1 - \alpha)) = \alpha$.*

Theorem 3 shows that we can estimate $c_n^*(1 - \alpha)$ through simulation: let $c_n^B(1 - \alpha)$ denote the $\alpha \cdot 100\%$ percentile of the test statistics $v(t_n^{*(1)}(\gamma)), \dots, v(t_n^{*(B)}(\gamma))$ obtained from B bootstrap samples. As $B \rightarrow \infty$, $c_n^B(1 - \alpha)$ becomes arbitrarily close to $c_n^*(1 - \alpha)$.

The above framework can also be used to derive a test based on $\inf t_n$. In this case the alternative hypothesis is that the benchmark model is worse than the competing model uniformly on Γ , i.e., under the alternative hypothesis the competing model dominates the benchmark uniformly. Below we consider a framework for forecast dominance based on the employment of two one-sided tests.

3.2 Feasible implementation of the tests

When $\Gamma \in \mathbb{R}^d$ contains infinitely many elements, as in our main examples, it is not possible to evaluate the test statistic in equation (13). Here we provide two numerical approximations for which Theorems 1 and 2 remain valid.

We first consider discretizations of Γ that become increasingly dense. Consider a grid of Γ with K_n elements Γ_n^i , such that $\sup_{\gamma, \gamma' \in \Gamma_n^i} |\gamma - \gamma'| < \delta_n$, for all $i = 1, \dots, K_n$, and let $\gamma_{n,i}$ be some point in Γ_n^i . We have the following approximation to our test statistic:

$$\widehat{\sup t_n} \equiv \max_{i=1, \dots, K_n} t_n(\gamma_{n,i}). \quad (15)$$

If Γ is a hyperrectangle in \mathbb{R}^d , a particularly convenient choice of K_n and $\{\Gamma_n^i\}_{i=1}^{K_n}$ derives from partitioning each dimension of Γ in v_n equal parts, which results in $K_n = v_n^d$.

The condition that $\delta_n \rightarrow 0$ implies that K_n grows quickly with large d . As a result, the calculation of $\widehat{\sup t_n}$ becomes problematic for large d . In such cases one can instead use Monte Carlo draws from some distribution J with uniformly positive density on Γ (e.g., a Uniform distribution) to obtain an approximation of the test statistic. Consider S_n independent draws

$\gamma^{(i)}$ from J , $i = 1, \dots, S_n$, and the approximation:

$$\widehat{\sup t_n} \equiv \max_{i=1, \dots, S_n} t_n(\gamma^{(i)}). \quad (16)$$

Proposition 1. *Let the assumptions of Theorem 1 hold. For some $K_n \rightarrow \infty$, such that $\delta_n \rightarrow 0$, as $n \rightarrow \infty$, $\widehat{v(t_n)} \xrightarrow{p} v(t_n)$. Moreover, $\widehat{v(t_n)} \xrightarrow{p} v(t_n)$ for $S_n \rightarrow \infty$ as $n \rightarrow \infty$.*

3.3 Detecting forecast dominance

A rejection of the null hypothesis H_0 in equation (2) constitutes evidence against model \mathcal{A} in favor of model \mathcal{B} at some values of γ in Γ . However, it is possible that model \mathcal{A} is superior to model \mathcal{B} for other points in Γ , and thus that neither model uniformly dominates the other.

To resolve this ambiguity we propose a framework using two one-sided tests: the first being the test of H_0 above, and the second being a test of the *opposite* null hypothesis:

$$H_0'' : E[-L_{t+1}(\gamma)] \leq 0 \quad \forall \gamma \in \Gamma$$

and which is implemented simply by using test statistic $\sup -t_n$.

Similar to forecast encompassing tests, see Chong and Hendry (1986), employing both tests results in one of the following four outcomes:

1. Fail to reject H_0 , reject H_0'' . \mathcal{A} significantly beats \mathcal{B} for some values of γ , and is not significantly beaten by \mathcal{B} for any γ . Thus \mathcal{A} dominates \mathcal{B} .
2. Reject H_0 , fail to reject H_0'' : Similar to outcome 1, but \mathcal{B} dominates \mathcal{A} .
3. Fail to reject both H_0 and H_0'' . Neither model significantly beats the other for any value of γ . Thus the models have statistically equal performance across γ .
4. Reject both H_0 and H_0'' : There are values of γ for which \mathcal{A} significantly beats \mathcal{B} , and values for which \mathcal{B} significantly beats \mathcal{A} . Thus there is no ordering of the models across all γ .

Outcomes 1 and 2 clearly reveal a preferred model. Outcome 3 indicates a lack of power to distinguish between the competing models (or actual equality of forecast performance across γ). Outcome 4 reveals an important sensitivity in the ranking of the models to the choice of shape parameter. Given that the above procedure involves two tests, each with non-zero Type I error

probability, and each with, inevitably, imperfect finite-sample power, we interpret the outcome of the above procedure as merely informative about forecast dominance, not as a formal test.

4 Simulation studies

In this section we evaluate the finite-sample performance of our proposed test in the three applications described in Section 2. In Section 4.1 we present an application for i.i.d. data, in Sections 4.2–4.3 we present two time series simulation designs, and in Section 4.4 we present the results.

4.1 Comparisons based on expected utility

We consider the difference in expected utility of two commonly-used portfolio management strategies: the equal-weighted portfolio and the minimum-variance portfolio. For an in-depth analysis of these and other portfolio strategies see DeMiguel et al. (2007). These strategies can be defined in terms of portfolio weight vectors:

$$\begin{aligned} w_t^{\text{eq}} &= \frac{1}{N} \iota, \\ w_t^{\text{mv}} &= \frac{1}{\iota' \Sigma_t^{-1} \iota} \Sigma_t^{-1} \iota, \end{aligned} \tag{17}$$

where Y_t is an $(N \times 1)$ vector of monthly excess returns with conditional covariance matrix Σ_t , and ι is an $(N \times 1)$ vector of ones. Denote the portfolio returns as $Y_t^{\text{eq}} = w_t^{\text{eq}'} Y_t$, and $Y_t^{\text{mv}} = w_t^{\text{mv}'} Y_t$. The feasible counterpart of w_t^{mv} depends on an estimate of Σ_t . In our simulation study we consider a simple rolling window estimate of the covariance matrix based on the most recent 120 observations, corresponding to 10 years of monthly data or six months of daily data.

We test whether the equal weighted and minimum-variance portfolio returns have equivalent expected utility across a range of levels of risk aversion. We model utility using the exponential utility function $u(y; \gamma) = -\exp\{-\gamma y\}/\gamma$. A wide range of values of the risk aversion parameter have been reported in the literature, ranging from near zero to as high as 60, see Bliss and Panigirtzoglou (2004). These authors estimate this parameter as being between 2.98 to 10.56, while DeMiguel et al. (2007) perform comparisons for investors with risk aversion ranging between 1 and 10. Based on this, we test for equal expected utility over $\Gamma = [1, 10]$.

We simulate excess returns Y_t according to a one-factor model, based on the DGP in the simulation study of DeMiguel et al. (2007). Let $Y_t = (Y_t^f, Y_{1,t}, \dots, Y_{N-1,t})'$, where Y_t^f denotes

the excess return on the factor portfolio and $Y_{i,t}$ denotes the $N - 1$ excess returns, generated as:

$$\begin{aligned} Y_{i,t} &= \alpha_i + \beta_i Y_t^f + \eta_{i,t}, \\ \nu_{i,t} &\sim \text{iid } \mathcal{N}(0, \sigma_{\eta,i}^2), \\ Y_t^f &\sim \text{iid } \mathcal{N}(\mu_f, \sigma_f^2). \end{aligned} \tag{18}$$

We follow the parameterization of DeMiguel et al. (2007), which resembles estimates that are commonly found in empirical studies. We set $\alpha_i = 0$, and $\beta_i = 0.5 + (i - 1)/(N - 1)$, for all $i = 1, \dots, N - 1$. Moreover, we set $\mu_f = 8\%$, and $\sigma_f = 16\%$. Finally, we let the idiosyncratic volatilities vary between 10% and 30%. However, unlike DeMiguel et al. (2007), who draw from the uniform distribution on [10%, 30%], we opt for deterministic cross-sectional variation between 10% and 30% by setting $\sigma_{\eta,i} = 10\% + 20\% \cdot \sin(\pi(i - 1)/(N - 1))$. We do so to facilitate the approximation of $E[L_t(\gamma)]$, which is required in the size experiment.

Given the portfolio strategies we consider, it is not generally possible to find a parameterization that implies $E[L_t(\gamma)] = 0$ for all $\gamma \in \Gamma$, which is the point in the null hypothesis (equation 2) least favorable to the alternative. In the size experiment we therefore test the null hypothesis $E[L_t(\gamma)] = \zeta_m(\gamma)$ instead of zero, where $\zeta_m(\gamma) \neq 0$ is the population value of $E[L_t(\gamma)]$, which we estimate using 100,000 simulations. (Note that $L_t(\gamma)$ depends on the rolling window sample covariance matrix, and this is used when computing its population expectation.) The power experiment tests the hypothesis $E[L_t(\gamma)] = 0$ for all $\gamma \in \Gamma$.

4.2 Forecast comparison via tail quantile forecasts of portfolio returns

We next study a scenario comparing quantile forecasts of portfolio returns implied by multivariate forecasting models. We simulate returns data using a GARCH-DCC model (Engle, 2002) with normal errors, parameterized to resemble the properties of daily asset returns.

We compare two widely-used models: (i) a GARCH-DCC model with normal errors, and (ii) a multivariate normal distribution with the RiskMetrics covariance estimator (Riskmetrics,

1996). We let the $N \times 1$ return vector Y_{t+1} follow a GARCH-DCC process

$$\begin{aligned}
Y_{t+1} &= \mu_{t+1} + H_{t+1}^{1/2} C_{t+1}^{1/2} \nu_{t+1}, \\
\nu_{t+1} &= (\nu_{t+1,1}, \dots, \nu_{t+1,N})' \sim iid N(0, I), \\
H_{t+1} &= \text{diag}(h_{t+1,1}, \dots, h_{t+1,N}), \\
h_{t+1,i} &= \omega_0 + \omega_1 h_{t,i} + \omega_2 h_{t,i} \nu_{t,i}^2, \\
C_{t+1} &= \text{diag}(\tilde{C}_{t+1})^{-1/2} \tilde{C}_{t+1} \text{diag}(\tilde{C}_{t+1})^{-1/2}, \\
\tilde{C}_{t+1} &= (1 - \xi_1 - \xi_2) \bar{C} + \xi_1 \tilde{C}_t + \xi_2 \text{diag}(\tilde{C}_t)^{1/2} \nu_t \nu_t' \text{diag}(\tilde{C}_t)^{1/2}, \\
\bar{C} &= [\bar{C}]_{ij}, \text{ where } [\bar{C}]_{ij} = 1 - \frac{|i-j|}{N}.
\end{aligned} \tag{19}$$

We choose GARCH parameters $\omega_0 = 0.05$, $\omega_1 = 0.10$, $\omega_2 = 0.85$ and DCC parameters $\xi_1 = 0.025$, $\xi_2 = 0.95$, to match time-varying volatility and correlation patterns commonly found in daily equity returns. We set $\mu_{t+1} = 0$ for simplicity. To fix the value of \bar{C} we use the covariance matrix generated by a Bartlett kernel with bandwidth set to N . This specification generates a diverse set of correlations, and ensures positive definiteness of \bar{C} .

We are interested in one-period-ahead α -quantile forecasts for portfolio returns $\tilde{Y}_{t+1}(\gamma) = \gamma' Y_{t+1}$, with $\alpha = 5\%$, for a range of portfolio weights $\gamma \in \Gamma$. The first forecast we consider is the optimal forecast, based on the GARCH-DCC model above. This forecast is given by

$$Q_{t,\alpha}^{\text{DCC}}(\gamma) = \Phi^{-1}(\alpha) \cdot \sqrt{\gamma' H_{t+1}^{1/2} C_{t+1} H_{t+1}^{1/2} \gamma}, \tag{20}$$

where $\Phi^{-1}(\alpha)$ denotes the α -quantile of the standard normal distribution.

The RiskMetrics forecast is

$$Q_{t,\alpha}^{\text{RM}}(\gamma) = \Phi^{-1}(\alpha) \cdot \sqrt{\gamma' \hat{\Sigma}_{t+1} \gamma}, \tag{21}$$

$$\text{where } \hat{\Sigma}_{t+1} = c_{\lambda,m} \sum_{j=0}^{m-1} \lambda^j (Y_{t-j} - \hat{\mu}_t) (Y_{t-j} - \hat{\mu}_t)' \tag{22}$$

where $\hat{\mu}_{t+1} = \frac{1}{m} \sum_{j=0}^{m-1} Y_{t-j}$ and $c_{\lambda,m}$ is a constant that normalizes the summed weights $\sum_{j=1}^m \lambda^j$ to one. As is standard for the RiskMetrics approach using daily returns, we set $\lambda = 0.94$.

We obtain the scores (losses) $S_{t+1}^i(\gamma)$ using the tick loss function, which is a consistent loss

function for the quantile, and is defined as

$$S_{t+1}^i(\gamma, \alpha) = (\mathbf{1}\{\tilde{Y}_{t+1}^p(\gamma) < Q_{t,\alpha}^i(\gamma)\} - \alpha)(Q_{t,\alpha}^i(\gamma) - \tilde{Y}_{t+1}^p(\gamma)), \text{ for } i \in \{\text{DCC}, \text{RM}\}, \quad (23)$$

and obtain loss differences $L_{t+1}(\gamma) = S_{t+1}^{\text{DCC}}(\gamma) - S_{t+1}^{\text{RM}}(\gamma)$.

We consider all portfolio return vectors with positive weights summing to one, i.e. $\Gamma = \{\gamma : \gamma_i \geq 0, i = 1, \dots, N, \sum_{i=1}^N \gamma_i = 1\}$. Γ is thus the $(N - 1)$ -simplex, and drawing uniformly from Γ is particularly easy using the Dirichlet distribution of order N with concentration parameters set to one. See Kotz et al. (2000, Ch. 49) for an elaborate treatment of the Dirichlet distribution.

As in the previous example, to study the finite-sample size of this test we consider the null hypothesis that $E[L_{t+1}(\gamma)] = \zeta_m(\gamma)$ instead of zero, where $\zeta_m(\gamma) \neq 0$ is the population value of $E[L_{t+1}(\gamma)]$, which we estimate using 100,000 simulations. The power experiment tests the null hypothesis $E[L_{t+1}(\gamma)] = 0$ for all $\gamma \in \Gamma$.

4.3 Forecast comparison via Murphy diagrams of quantile forecasts

Under mild regularity conditions (see, e.g., Gneiting, 2011a) the α -quantile of a random variable Y_t is elicitable using the ‘‘generalized piecewise linear’’ (GPL) class of scoring rules:

$$S(Y_{t+1}, x; \alpha, g) = (\mathbf{1}(Y_{t+1} < x) - \alpha)(g(x) - g(Y_{t+1})), \quad (24)$$

where $g(\cdot)$ is a non-decreasing function. A commonly used member of this family is the tick loss function, which sets $g(z) = z$, and which was used in the previous section. In this example we test for differences in predictive ability of competing α -quantile forecasts across the set of *all* consistent scoring rules for α -quantiles, denoted \mathcal{S}_{GPL}^α , using the mixture representation for this class of loss functions presented in Ehm et al. (2016):

$$S(Y_{t+1}, x; \alpha, g) = \int_{-\infty}^{\infty} \tilde{S}(Y_{t+1}, x; \alpha, \gamma) dH(\gamma; g) \quad (25)$$

$$\text{where } \tilde{S}(Y_{t+1}, x; \alpha, \gamma) = (\mathbf{1}(Y_{t+1} < x) - \alpha)(\mathbf{1}(\gamma < x) - \mathbf{1}(\gamma < Y_{t+1})), \quad (26)$$

where $H(\cdot; g)$ some non-negative measure.

Consider two forecasts $Q_{t,\alpha}^A$ and $Q_{t,\alpha}^B$. We say that $Q_{t,\alpha}^A$ uniformly outperforms $Q_{t,\alpha}^B$ if

$$E[L(Y_t, Q_{t,\alpha}^A, Q_{t,\alpha}^B, \alpha, g)] \equiv E[S(Y_{t+1}, Q_{t,\alpha}^A, \alpha, g) - S(Y_{t+1}, Q_{t,\alpha}^B, \alpha, g)] \leq 0 \quad \forall S \in \mathcal{S}_{GPL}^\alpha \quad (27)$$

Corollary 1 in Ehm et al. (2016) establishes that equation (27) is implied by

$$E[\tilde{L}(Y_{t+1}, Q_{t,\alpha}^A, Q_{t,\alpha}^B, \alpha, \gamma)] \equiv E[\tilde{S}(Y_{t+1}, Q_{t,\alpha}^A, \alpha, \gamma) - \tilde{S}(Y_{t+1}, Q_{t,\alpha}^B, \alpha, \gamma)] \leq 0 \quad \forall \gamma \in \mathbb{R} \quad (28)$$

We can therefore test for superior predictive ability by testing the above condition. It should be noted that our theory requires the range of γ , denoted Γ , to be bounded, and so we cannot test over all $\gamma \in \mathbb{R}$. However, in practice we can make Γ large enough to cover all relevant parameter values, since $\tilde{S}(Y_{t+1}, Q_{t,\alpha}^A; \alpha, \gamma) - \tilde{S}(Y_{t+1}, Q_{t,\alpha}^B; \alpha, \gamma)$ is known to be identically zero for $\gamma \notin [\min(Y_{t+1}, Q_{t,\alpha}^A, Q_{t,\alpha}^B), \max(Y_{t+1}, Q_{t,\alpha}^A, Q_{t,\alpha}^B)]$.

In small samples it can occur that $\tilde{S}(Y_{t+1}, Q_{t,\alpha}^A; \alpha, \gamma) - \tilde{S}(Y_{t+1}, Q_{t,\alpha}^B; \alpha, \gamma) = 0$ for all observations in a given sample, for some values of γ . As a result, $\hat{\sigma}_n^2(\gamma) = 0$ for these values of γ . To circumvent this we fix $\sigma_n^2(\gamma) = 1$ for all $\gamma \in \Gamma$, i.e. we consider a test based on $t_n(\gamma) = \sqrt{n}\bar{L}_n(\gamma)$ instead of $t_n(\gamma) = \sqrt{n}\bar{L}_n(\gamma)/\hat{\sigma}_n(\gamma)$. The p -values remain valid under the bootstrap. The HAC covariance estimators used in the calculation of the multivariate Wald test and the Bonferroni-corrected test suffer from the same singularity. However, inference is no longer valid for these tests, because the limit law of these test statistics is no longer standard without studentization.

We use the same quantile forecast models as in the simulation design in the previous section, but set $N = 1$, so that the quantile forecasts defined in equations (20) and (21) are obtained from a GARCH model instead of a GARCH-DCC model. For an additional comparison, we also consider a rolling window sample quantile estimated over the previous 250 days.

As in the previous examples, to study the finite-sample size of this test we consider the null hypothesis that $E[L_{t+1}(\gamma)] = \zeta_m(\gamma)$ instead of zero, where $\zeta_m(\gamma) \neq 0$ is the population value of $E[L_{t+1}(\gamma)]$, which we estimate using 100,000 simulations. The power experiment tests the null hypothesis $E[L_{t+1}(\gamma)] = 0$ for all $\gamma \in \Gamma$.

4.4 Simulation results

Table 1 presents rejection rates for the size and power experiments introduced in Section 4.1, based on 1,000 Monte Carlo simulations. We consider out-of-sample period lengths $n = 500$ and 2,000 observations, two common sample sizes in applied work. We consider increasingly large grids of Γ , with the number of grid points set to $K_n = 1, 10, 50, 100$, and 250. We obtain critical values using $B = 1,000$ bootstrap samples, and for the block length l we follow a standard rule-of-thumb in the HAC literature and set $l = 4(n/100)^{2/9}$, which satisfies our Assumption 5.

Table 1: Small sample rejection rates of superior expected utility tests on equal weighted and minimum-variance portfolio strategies

K_n	Panel A: Size				Panel B: Power			
	n=500		n=2,000		n=500		n=2,000	
	Bonf.	sup- t	Bonf.	sup- t	Bonf.	sup- t	Bonf.	sup- t
1	0.08	0.09	0.02	0.02	0.00	0.00	0.00	0.00
10	0.02	0.10	0.02	0.10	0.99	1.00	1.00	1.00
50	0.02	0.10	0.01	0.04	0.99	1.00	1.00	1.00
100	0.02	0.09	0.01	0.07	0.97	1.00	1.00	1.00
250	0.01	0.13	0.00	0.07	0.95	1.00	1.00	1.00

Note: This table presents the rejection rates of the proposed one-sided test (sup- t), as well as the benchmark test (Bonferroni). The data is generated according to equation (18), and the equal-weighted and minimum-variance portfolio strategies are given in Equation (17). The minimum-variance portfolio weights are estimated using a rolling window of $m = 120$ observations. The out-of-sample period consists of $n = 500$, and 2,000 observations. We consider discrete grids of $\Gamma = [1, 10]$ formed using $K_n = 1, 10, 50, 100$, and 250 equally spaced grid points.

For comparison with the proposed test, we consider applying the familiar Bonferroni multiple-comparison correction. It should be noted that this correction can only be applied at a finite number of points, M , in Γ , and therefore cannot generally test over all Γ . A one-sided α -level test using the Bonferroni correction rejects the null hypothesis if, for at least one $\gamma \in \Gamma_M$, we find $n\bar{L}_n(\gamma)/\tilde{\sigma}_n(\gamma) > z_{1-\alpha/M}^{-1}$, with $z_{1-\alpha/M}^{-1}$ denoting the $(1 - \alpha/M)$ -quantile of the standard normal distribution. As usual for Bonferroni-corrected tests, the critical value is much larger than for the individual tests ($z_{1-\alpha/M}^{-1}$ rather than $z_{1-\alpha}^{-1}$) which can lead to low power.

From Panel A of Table 1 we observe that the sup- t test is somewhat oversized for $n = 500$ but approximately correctly-sized for $n = 2,000$. Reassuringly, we observe that the sup- t tests are stable in terms of rejection rates once K_n is moderately large, indicating robustness to this tuning parameter. The Bonferroni-based test has satisfactory size control for the smaller sample size, and becomes conservative for the larger sample size. Panel B shows rejection rates in the power experiment, and we observe the sup- t has good power for both sample sizes.

Table 2 presents small sample rejection rates of the size and power experiments introduced in Section 4.2, for portfolios composed of 30 assets. Results are presented for sample sizes of $n = 500$ and 2,000 observations, and six sets of weight vectors, which are drawn as follows. We first consider a set of 31 deterministic weight vectors: the equal weighted portfolio weights

Table 2: Small sample rejection rates of quantile forecast tests, for differences between multivariate GARCH-DCC and RiskMetrics models

S_n	Panel A: Size				Panel B: Power			
	n=500		n=2,000		n=500		n=2,000	
	Bonf.	sup- t	Bonf.	sup- t	Bonf.	sup- t	Bonf.	sup- t
31	0.01	0.01	0.01	0.01	0.16	0.15	0.62	0.65
50	0.01	0.01	0.00	0.01	0.11	0.14	0.55	0.64
100	0.00	0.01	0.00	0.01	0.07	0.14	0.43	0.63
250	0.00	0.01	0.00	0.01	0.04	0.14	0.28	0.64
500	0.00	0.01	0.00	0.01	0.03	0.14	0.21	0.64
1000	0.00	0.01	0.00	0.01	0.02	0.13	0.15	0.64

Note: This table presents the rejection rates of the proposed one-sided test (sup- t), as well as the benchmark Bonferroni test. The quantile forecasts for the portfolio returns from the GARCH-DCC and multivariate RiskMetrics models are defined in equations (20) and (21). The data is generated as in equation (19) with $N = 30$. We test at 31 fixed portfolio weight vector being the equal weighted portfolio vector and the 30 basis vectors, as well as $S_n - 31$ weight vectors drawn uniformly from the unit simplex.

and the 30 basis vectors. Subsequently we randomly draw $S_n - 31$ weight vectors from J , for $S_n = 50, 100, 250, 500$ and 1,000. We use $B = 1,000$ bootstrap samples to obtain critical values.

Panel A of Table 2 provides rejection rates in the size experiment. We observe that the sup- t test is conservative, though less than the benchmark Bonferroni test. Panel B shows rejection rates in the power experiment, and as in the previous section, we observe the sup- t test rejection probabilities are stable across values of S_n . Power is greater than the Bonferroni-based test, and that test's power declines monotonically as S_n increases.

Finally, Table 3 provides small sample rejection rates of the tests in size and power experiments introduced in Section 4.3. The columns labeled "RM" compare GARCH forecasts with RiskMetrics forecasts, and the columns labeled "RW" compares GARCH forecasts with rolling window forecasts. We again consider sample sizes of $n = 500$ and $n = 2,000$. We consider grids of Γ with $K_n = 50, 100,$ and 250 grid points equally spaced over the interval $[-20, 0]$. We select this interval because outside of it the elementary score differences are equal to zero in almost all realizations. As we cannot always (across values of γ in the elementary scores) compute the asymptotic covariance required to obtain the benchmark Bonferroni-corrected tests, we do not implement those here. Instead, we present the results of a standard Diebold-Mariano test based on the tick loss function as benchmark. These are reported in the rows labeled "1".

Table 3: Small sample rejection rates of Murphy diagram tests, for quantile forecast differences between GARCH, RiskMetrics, and Rolling Window models

		Panel A: Size				Panel B: Power			
		n=500		n=2,000		n=500		n=2,000	
Test	K_n	RM	RW	RM	RW	RM	RW	RM	RW
DM	1	0.04	0.05	0.04	0.06	0.19	0.45	0.49	0.95
sup- t	50	0.03	0.07	0.03	0.05	0.06	0.36	0.14	0.81
sup- t	100	0.02	0.09	0.04	0.08	0.06	0.39	0.13	0.85
sup- t	250	0.01	0.08	0.04	0.06	0.04	0.36	0.15	0.88

Note: This table presents the rejection rates of the proposed one-sided test (sup- t) comparing forecasts from a GARCH model with those from the RiskMetrics (RM) and Rolling Window (RW) forecasts. The Diebold-Mariano test (DM) using the tick loss function is also presented. The quantile forecasts from the GARCH and RiskMetrics models are given in equations (20) and (21), with $N = 1$. We consider out-of-sample period lengths $n = 500$, and 2,000, and discrete grids of $\Gamma = [-20, 0]$ with $K_n = 50, 100$, and 250 equally-spaced points.

Panel A of Table 3 provides rejection rates for the size experiments. We observe that the rejection rates are close to their nominal level for both sample sizes, though the sup- t test is somewhat conservative at $n = 500$. The benchmark Diebold-Mariano test using the tick loss function, is also approximately correctly sized. The rejection rates of the sup- t test are stable across values K_n , indicating robustness to this tuning parameter. Panel B of Table 3 provides rejection rates in the power experiment. In this design, the RiskMetrics and GARCH forecasts generate losses that are not very different and so power is low in that case, especially for the smaller sample size. At $n = 2,000$ the proposed test has power against the alternative, but the benchmark Diebold-Mariano (DM) test turns out to be more powerful in this design. It should of course be noted that the DM test considers a different null to ours; that test considers only a single loss function rather than a continuum of loss functions.

5 Applications to equity return forecast environments

5.1 Comparing the utility of two portfolio strategies

In this section we compare equal weighted and minimum-variance portfolio strategies in terms of average utility, across a range of levels of risk aversion. We use monthly returns data on 30

Table 4: p -values from tests of superior expected utility from equal-weighted and minimum-variance portfolio strategies

K_n	Risk aversion in [1,10]		Risk aversion in [1,5]		Risk aversion in [5,10]	
	Bonf.	sup- t	Bonf.	sup- t	Bonf.	sup- t
Panel A: H_1: EW superior for some γ						
1	0.02	0.01	0.02	0.02	0.56	0.57
10	0.02	0.00	0.02	0.00	1.00	0.22
50	0.11	0.00	0.11	0.00	1.00	0.17
100	0.22	0.00	0.22	0.00	1.00	0.22
250	0.56	0.01	0.56	0.01	1.00	0.20
Panel B: H_1: MV superior for some γ						
1	0.97	0.97	0.98	0.98	0.44	0.43
10	1.00	0.18	1.00	0.86	1.00	0.19
50	1.00	0.20	1.00	0.85	1.00	0.22
100	1.00	0.24	1.00	0.83	1.00	0.17
250	1.00	0.21	1.00	0.85	1.00	0.20

Note: This table presents p -values from tests for superior predictive accuracy. The equal-weighted and minimum-variance portfolio strategies are given in equation (17). The data consists of monthly returns of 30 industry portfolios and runs from September 1926 to December 2017. The three panels consider three ranges of values for the risk aversion parameter. K_n indicates the number of (equally-spaced) grid points.

U.S. industry portfolios, from September 1926 to December 2017, a total of 1,098 observations.¹ The minimum-variance portfolio weights are estimated over a rolling window of $m = 120$ months. We use the exponential utility function $u(y; \gamma) = -\exp\{-\gamma y\}/\gamma$, with $\gamma \in [1, 10]$, which covers all risk aversion parameter values considered in DeMiguel et al. (2007).

Table 4 provides p -values for the proposed new test, as well as the benchmark Bonferroni-adjusted tests. The rows labeled $K_n = 1$ tests the hypothesis of equal predictive ability for one particular value of risk aversion $\gamma = 2.5$, while the remaining rows use $K_n = 10, 50, 100, 250$ equally-spaced points in $[1, 10]$. We set $B = 1,000$.

In the left panel Table 4 we consider $\Gamma = [1, 10]$. When $K_n = 1$, the sup- t test rejects the null of weakly greater utility from the minimum variance portfolio compared with the equal-weighted portfolio and fails to reject the opposite null, indicating that the equal-weighted

¹The returns can be obtained from the data library at http://mba.tuck.dartmouth.edu/pages/faculty/ken.french/data_library.html

strategy dominates the minimum variance portfolio for this particular risk aversion. When we increase K_n , and test the hypothesis of equal average utility over the entirety of Γ , we find the much stronger result that the equal-weighted portfolio dominates the minimum variance portfolio across all values for the risk aversion parameter. In contrast, the Bonferroni-corrected test fails to reject the null hypothesis for larger values of K_n . This finding is consistent with the conservativeness of this test as K_n increases documented in the simulation study in the previous section.

Figure 1 shows the sample mean of expected utility differences and the pointwise 95% confidence bounds for each $\gamma \in [1, 10]$. We observe that for γ less than about 3, the confidence interval does not include zero, indicating that for lower levels of risk aversion the equal-weighted portfolio significantly outperforms the minimum variance portfolio. For higher levels of risk aversion, the pointwise confidence intervals either contain zero, or lie below zero.

The middle and right panels of Table 4 show results for $\Gamma = [1, 5]$ and $\Gamma = [5, 10]$, to examine the sensitivity of the conclusions of these tests to the range of values of risk aversion considered. Consistent with Figure 1, for less risk averse investors, with $\Gamma = [1, 5]$, we observe that the equal-weighted strategy dominates the minimum variance strategy, while we cannot distinguish between the strategies over $\Gamma = [5, 10]$.

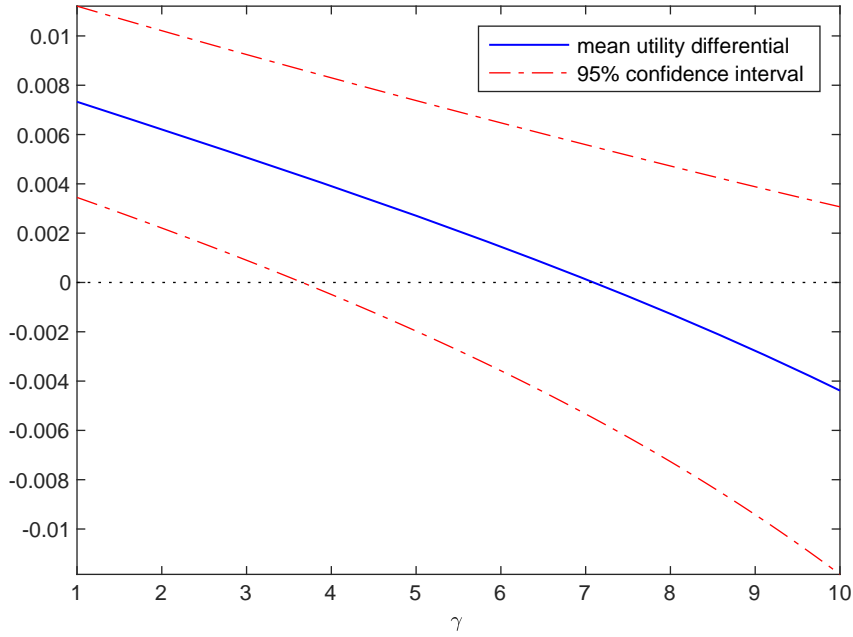
5.2 Quantile forecasts of portfolio returns from multivariate models

In this analysis we compare two multivariate volatility models, a GARCH-DCC model and the RiskMetrics model, by the quality of their forecasts for the 5% Value-at-Risk (i.e., the 5% quantile) of the returns on portfolios of the underlying assets. We use daily returns on the same 30 U.S. industry portfolios as in the previous section, over the period January 1998 to December 2017, a total of 5,032 observations. The GARCH-DCC model is estimated over a rolling window of 1,000 observations. We set $B = 1,000$.

We consider the following sets of portfolios. For $S_n = 1$ we consider only the equal weighted portfolio. For $S_n = 31$ we consider the equal weighted portfolio and all 30 single-asset portfolios. For $S_n > 31$, we additionally consider $S_n - 31$ random weight vectors drawn from the 30-dimensional simplex.

Table 5 presents p -values of the tests of equal or superior predictive ability. We show results for the full sample, as well as the first and second halves of the sample period. The results for the full sample show that the models have approximately equal performance, since the test

Figure 1: Utility differential of equal weighted and minimum-variance portfolio strategies



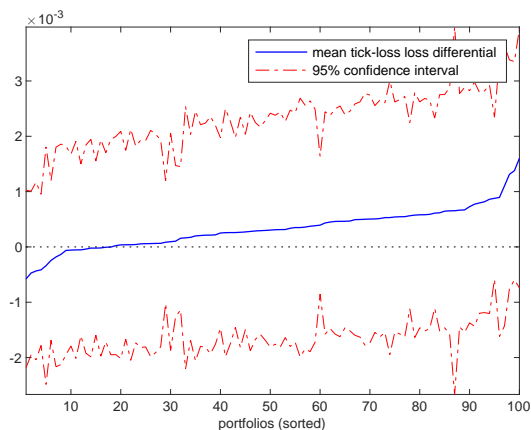
Note: This figure plots the the sample mean of utility from the equal-weighted portfolio strategy *minus* the sample mean of utility derived from the minimum-variance portfolio strategies. The strategies are given equation (17). The data consists of monthly returns of 30 industry portfolios from September 1926 to December 2017. We consider risk aversion, γ , in the range $[1, 10]$.

cannot reject the null in favor of either RiskMetrics or DCC; the p -values are all well above 0.05.

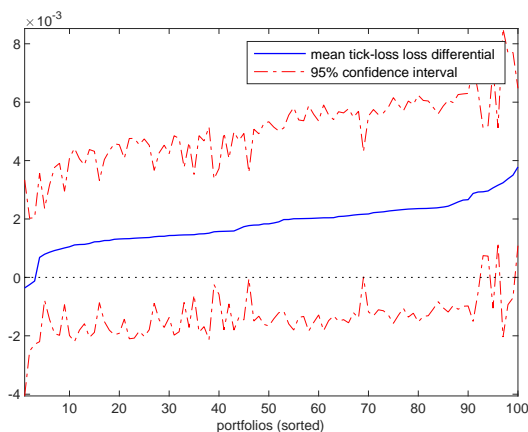
In the middle panel of Table 5, which examines performance in the first half of our sample period, we find that we can reject in favor of the RiskMetrics model, but not the reverse, indicating that the RiskMetrics model dominates the GARCH-DCC model in this sub-sample. In the second sub-sample, similar to the full sample, we cannot reject the hypotheses of superior predictive ability, indicating an inability to distinguish between these forecasting models.

Figure 2 plots the sample mean of the tick loss differences and the 95% confidence bounds for the first 100 portfolio weights that we draw, over the full sample and subsamples, and sorted on mean tick loss difference. In the first half of the sample we indeed find that the RiskMetrics forecasts perform better, although for only few portfolio vectors do we observe (pointwise) confidence intervals that exclude zero. In the second half of the sample the average tick loss differences are generally negative, but not significantly different from zero.

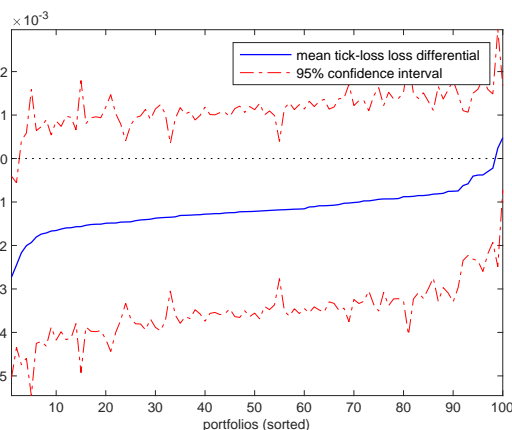
Figure 2: Tick-loss differential for tail quantile forecasts generated by the multivariate GARCH-DCC and RiskMetrics models



(a) Full sample



(b) First sub-sample



(c) Second sub-sample

Note: This figure plots the sample average of the tick loss from the GARCH-DCC forecasts *minus* the sample mean tick loss of the RiskMetrics forecasts. The forecasts are defined in equations (20) and (21). The mean tick loss differences are shown for $S_n = 100$ portfolio weight vectors, with 31 of these being the equal weighted portfolio vector and the 30 single-asset portfolio vectors, and the rest drawn uniformly from the 30-dimensional simplex. The portfolios are presented sorted on mean tick loss difference. The data consists of daily returns for 30 Industry portfolios and runs from January 1998 to December 2017.

Table 5: p -values of quantile forecast tests, for differences between multivariate GARCH-DCC and RiskMetrics models

S_n	Full sample		First sub-sample		Second sub-sample	
	Bonf.	sup- t	Bonf.	sup- t	Bonf.	sup- t
Panel A: H_1 : RiskMetrics superior for some γ						
1	0.39	0.39	0.15	0.15	0.84	0.85
31	1.00	0.62	0.02	0.02	1.00	0.87
100	1.00	0.64	0.08	0.03	1.00	0.90
250	1.00	0.59	0.20	0.03	1.00	0.91
500	1.00	0.54	0.40	0.04	1.00	0.90
1000	1.00	0.58	0.80	0.03	1.00	0.91
Panel B: H_1 : DCC superior for some γ						
1	0.61	0.56	0.85	0.85	0.16	0.11
31	1.00	0.93	1.00	0.99	0.17	0.08
100	1.00	0.93	1.00	1.00	0.56	0.09
250	1.00	0.93	1.00	1.00	1.00	0.09
500	1.00	0.92	1.00	1.00	1.00	0.09
1000	1.00	0.93	1.00	0.99	1.00	0.11

This table presents p -values from tests for superior predictive accuracy. The quantile forecasts for the portfolio returns from the GARCH-DCC and multivariate RiskMetrics models are defined in equations (20) and (21). The data consists of twenty years of daily returns of 30 industry portfolios from January 1998 to December 2017. We consider S_n portfolio weight vectors. When $S_n = 1$ we use equal weights. For $S_n \geq 31$ we test at 31 fixed portfolio weight vectors (the equal-weighted portfolio vector and the 30 basis vectors) and $S_n - 31$ weight vectors drawn uniformly from the 30-dimensional unit simplex.

5.3 Quantile forecast comparison via Murphy diagrams

Our final empirical analysis compares two forecast models for the 5% quantile (i.e., the 5% Value-at-Risk) of a single asset. We compare three models: a GARCH(1,1) model (Bollerslev, 1986), the RiskMetrics model, and a simple rolling window sample quantile calculated over the previous 250 days. We use the same data as the previous sub-section: daily returns on 30 U.S. industry portfolios, over the period January 1998 to December 2017, a total of 5,032 observations. The GARCH model is estimated over a rolling window of 1,000 observations. The parameter space for the elementary scoring rule shape parameter (see Section 4.3 for details) is $\Gamma = [-20, 0]$, and we consider an increasingly fine grid of equally-spaced points in this space when implementing the new tests. We set $B = 5,000$.

We implement the tests for each of the 30 industry portfolio returns separately. In Table 6 we present detailed results for a single representative portfolio (the “Transportation” industry portfolio), and in Table 7 we present a summary of the results across all 30 industry portfolios.

Panel A of Table 6 compares RiskMetrics and GARCH forecasts. We find that the benchmark Diebold and Mariano (1995) test using the tick loss function fails to reject both nulls, with a p -values well above 0.05. The sup- t test rejects the null of weakly superior GARCH forecasts with p -values around 0.03, but does not reject in the opposite direction, indicating that the RiskMetrics forecasts dominate the GARCH forecasts for the Transportation industry portfolio.

The upper panel of Figure 3 shows the “Murphy diagram” for this comparison, applied to the Transportation portfolio, and reveals that for most values of the elementary scoring rule parameter the GARCH and RiskMetrics forecasts have similar average losses. For values of the parameter around -1 the GARCH forecast significantly outperforms the RiskMetrics forecast, with the pointwise confidence intervals being far from zero, whereas the RiskMetrics forecast shows some (pointwise) significant outperformance for values of the parameter around -2 . The results of the tests in Table 6, however, indicate that the RiskMetrics forecasts dominate the GARCH forecasts overall.

Panel B of Table 6 compares the GARCH forecast with the rolling window sample quantile forecast. The sup- t test finds no evidence that the rolling window sample quantile outperforms the GARCH forecast for any value of elementary scoring rule parameter, whereas we do find evidence in the opposite direction, indicating that the GARCH forecast dominates the rolling window forecasts for the Transportation industry portfolio. The lower panel of Figure 3 shows that the difference in average loss is negative almost everywhere, consistent with the tests in Table 6, and the pointwise confidence intervals exclude zero for a large part of the parameter space.

Panel A of Table 7 compares Riskmetrics and GARCH forecasts across 30 industry portfolios and reports the proportion of portfolios in each of the four “forecast dominance” outcomes based on the sup- t tests discussed in Section 3.3. Focusing on the $K_n = 1,000$ row, we see that for 17% of the 30 portfolios either the RiskMetrics model dominates (10%) or the GARCH model dominates (7%), while for the remaining portfolios we are unable to statistically distinguish the performance of these two models. Consistent with this, results for the one-sided Diebold and Mariano (1995) tests using tick loss (not reported in the interests of space) do not reject for

Table 6: Quantile forecast comparison tests for the Transportation industry portfolio

K_n	H_1 : Comp. superior for some γ		H_1'' : GARCH superior for some γ	
	tick loss	sup- t	tick loss	sup- t
Panel A: RiskMetrics vs. GARCH				
1	0.38	-	0.62	-
50	-	0.01	-	0.06
100	-	0.04	-	0.21
250	-	0.05	-	0.25
500	-	0.03	-	0.24
1000	-	0.04	-	0.24
Panel B: Rolling window vs. GARCH				
1	1.00	-	0.00	-
50	-	0.96	-	0.00
100	-	0.98	-	0.00
250	-	0.99	-	0.00
500	-	0.98	-	0.00
1000	-	0.99	-	0.00

Note: This table presents p -values from tests comparing the predictive accuracy of forecasts of the 5% quantile for daily returns obtained from GARCH, RiskMetrics and rolling window models. We consider tests based on the elementary scoring rules for the quantile, as well as a test using the tick loss function, applied to daily returns on the Transportation industry portfolio. We consider equally-spaced discrete grids of $\Gamma = [-20, 0]$ with K_n grid points. The tick loss based tests use only a single loss function, and are reported in the rows labeled $K_n = 1$. The model referred to in the panel labeled “Comp. superior” is RiskMetrics in Panel A and rolling window in Panel B.

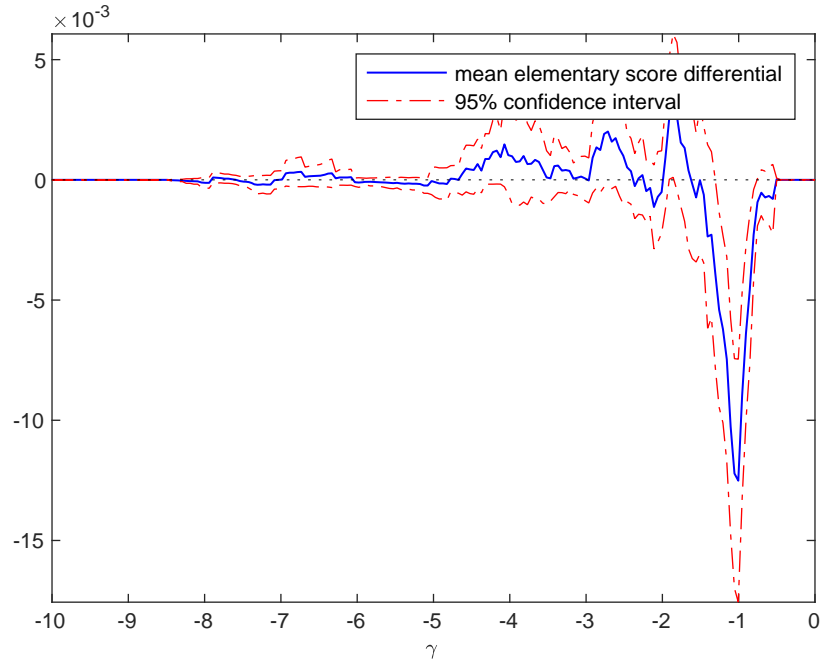
any portfolio in any direction.

In contrast, Panel B of Table 7 reveals a clear ordering of the rolling window and GARCH forecasts: the “forecast dominance” outcomes reveal that the GARCH model provides superior forecasts to the rolling window model, for all 30 portfolios. One-sided Diebold and Mariano (1995) tests (unreported) using only the tick loss function arrive at the same conclusion.

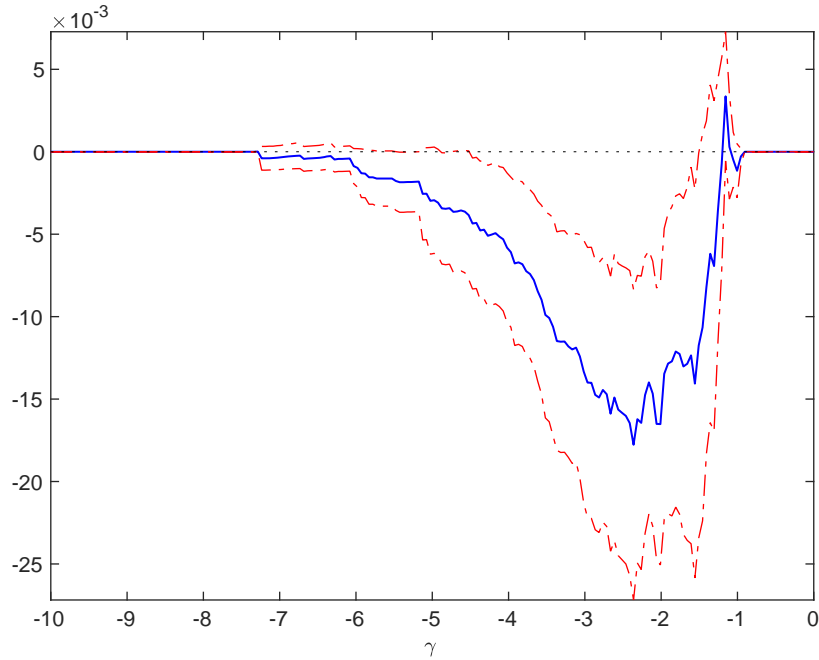
6 Concluding remarks

In many empirical applications, researchers are faced with the problem of comparing forecasts using a loss function that contains a shape parameter; examples include comparisons using average utility across a range of values for the level of risk aversion, and comparisons using

Figure 3: Murphy diagrams for quantile forecasts for the Transportation portfolio returns



(a) GARCH vs. RiskMetrics



(b) GARCH vs. rolling window sample quantile

Note: This figure plots the sample mean of the elementary scoring rule of the GARCH forecasts *minus* the sample mean of the elementary scoring rule of the RiskMetrics forecasts. The forecasts are 5% quantile forecasts for daily returns of the US Transportation industry index from January 1998 to December 2017. The elementary scoring rules are indexed by the scalar parameter $\gamma \in [-10, 0]$.

Table 7: Quantile forecast comparison test results across all 30 industry portfolios.

K_n	Forecast dominance outcomes			
	(1) Comp. dominates	(2) GARCH dominates	(3) No rejections	(4) No ordering
Panel A: RiskMetrics vs. GARCH				
50	0.17	0.00	0.83	0.00
100	0.23	0.10	0.67	0.00
250	0.10	0.03	0.83	0.03
500	0.10	0.07	0.80	0.03
1000	0.10	0.07	0.83	0.00
Panel B: Rolling window vs. GARCH				
50	0.00	1.00	0.00	0.00
100	0.00	1.00	0.00	0.00
250	0.00	1.00	0.00	0.00
500	0.00	1.00	0.00	0.00
1000	0.00	1.00	0.00	0.00

Note: This table presents the proportion of forecast dominance outcomes across 30 industry portfolios, using tests comparing the predictive accuracy of forecasts of the 5% quantile for daily returns obtained from GARCH, RiskMetrics and rolling window models. We consider equally-spaced discrete grids of $\Gamma = [-20, 0]$ with K_n grid points. The forecast dominance results in the final four columns are based on the decision rule described in Section 3.3. The model referred to in the column labeled “Comp. dominates” is RiskMetrics in Panel A and rolling window in Panel B.

characteristics of a portfolio return across a range of values for the portfolio weight vector. We propose new forecast comparison tests, in the spirit of Diebold and Mariano (1995) and Giacomini and White (2006), that may be applied in such applications. We consider tests for superior forecast accuracy across the entire range of values of the loss function parameter. The asymptotic properties of the test statistics are derived using block bootstrap theory for empirical processes, see Bühlmann (1995).

We show via an extensive simulation study that the tests have satisfactory finite sample properties, unlike the leading existing alternative which breaks down when a large number of values of the shape parameter is considered. We illustrate the new tests in three empirical applications: comparing portfolio strategies using average utility across a range of levels of risk aversion; comparing multivariate volatility models via their Value-at-Risk forecasts for portfolios

of the underlying assets across a range of values for the portfolio weight vector; and comparisons using recently-proposed “Murphy diagrams” (Ehm et al., 2016) for classes of consistent scoring rules for quantile forecasting.

This paper leaves open some interesting avenues for future research. If a functional CLT for triangular arrays of weakly dependent sequences can be invoked, then a time series analog of the approach in Andrews and Shi (2013) may be pursued, with the potential to improve power. In a different direction, our setting, where the researcher wishes to consider a range of values for the loss function parameter rather than a single value, may prompt one to consider a Bayesian approach with some prior distribution on the loss function parameter. We leave these explorations for subsequent research.

References

- Andrews, D. and Shi, X. (2013). Inference based on conditional moment inequalities. *Econometrica*, 81:609–666.
- Andrews, D. W. (1992). Generic Uniform Convergence. *Econometric Theory*, 8(2):241–257.
- Andrews, D. W. (1994). Empirical Process Methods in Econometrics. In Engle, R. F. and McFadden, D. L., editors, *Handbook of Econometrics, Volume 4*, pages 2247–2294. Elsevier.
- Bliss, R. R. and Panigirtzoglou, N. (2004). Option-Implied Risk Aversion Estimates. *Journal of Finance*, 59(1):407–446.
- Boussama, F., Fuchs, F., and Stelzer, R. (2011). Stationarity and Geometric Ergodicity of BEKK Multivariate GARCH Models. *Stochastic Processes and their Applications*, 121(10):2331–2360.
- Bradley, R. C. et al. (2005). Basic Properties of Strong Mixing Conditions. A Survey and Some Open Questions. *Probability Surveys*, 2:107–144.
- Bühlmann, P. (1995). The Blockwise Bootstrap for General Empirical Processes of Stationary Sequences. *Stochastic Processes and Their Applications*, 58(2):247–265.
- Chong, Y. Y. and Hendry, D. F. (1986). Econometric Evaluation of Linear Macroeconomic Models. *Review of Economic Studies*, 53:671–690.

- Clark, T. E. and McCracken, M. W. (2001). Tests of Equal Forecast Accuracy and Encompassing for Nested Models. *Journal of Econometrics*, 105(1):85–110.
- Davydov, Y., Lifshits, M. A., and Smorodina, N. (1998). *Local Properties of Distributions of Stochastic Functionals*. American Mathematical Society.
- DeMiguel, V., Garlappi, L., and Uppal, R. (2007). Optimal Versus equal weighted Diversification: How Inefficient is the 1/N Portfolio Strategy? *Review of Financial Studies*, 22(5):1915–1953.
- Diebold, F. X. and Mariano, R. S. (1995). Comparing Predictive Accuracy. *Journal of Business and Economic Statistics*, 13(3):253–263.
- Doukhan, P., Massart, P., and Rio, E. (1994). The Functional Central Limit Theorem for Strongly Mixing Processes. *Annales de l’IHP Probabilités et statistiques*, 30(1):63–82.
- Doukhan, P., Massart, P., and Rio, E. (1995). Invariance Principles for Absolutely Regular Empirical Processes. *Annales de l’I.H.P. Probabilités et Statistiques*, 31(2):393–427.
- Ehm, W., Gneiting, T., Jordan, A., and Krüger, F. (2016). Of Quantiles and Expectiles: Consistent Scoring Functions, Choquet Representations and Forecast Rankings. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 78(3):505–562.
- Engle, R. (2002). Dynamic Conditional Correlation: A Simple Class of Multivariate Generalized Autoregressive Conditional Heteroskedasticity Models. *Journal of Business and Economic Statistics*, 20(3):339–350.
- Engle, R. and Colacito, R. (2006). Testing and Valuing Dynamic Correlations for Asset Allocation. *Journal of Business and Economic Statistics*, 24(2):238–253.
- Fissler, T. and Ziegel, J. F. (2016). Higher Order Elicitability and Osband’s Principle. *Annals of Statistics*, 44(4):1680–1707.
- Fleming, J., Kirby, C., and Ostdiek, B. (2001). The Economic Value of Volatility Timing. *Journal of Finance*, 56(1):329–352.
- Giacomini, R. and White, H. (2006). Tests of Conditional Predictive Ability. *Econometrica*, 74(6):1545–1578.

- Gneiting, T. (2011a). Making and Evaluating Point Forecasts. *Journal of the American Statistical Association*, 106(494):746–762.
- Gneiting, T. (2011b). Quantiles as Optimal Point Forecasts. *International Journal of Forecasting*, 27(2):197–207.
- Hand, D. J. (1998). Data Mining: Statistics and More? *American Statistician*, 52(2):112–118.
- Hansen, P. R. (2005). A test for superior predictive ability. *Journal of Business and Economic Statistics*, 23(4):365–380.
- Hubrich, K. and West, K. D. (2010). Forecast evaluation of small nested model sets. *Journal of Applied Econometrics*, 25:574–594.
- Jin, S., Corradi, V., and Swanson, N. R. (2017). Robust Forecast Comparison. *Econometric Theory*, 33(6):1306–1351.
- Kole, E., Markwat, T., Opschoor, A., and Van Dijk, D. (2017). Forecasting Value-at-Risk under Temporal and Portfolio Aggregation. *Journal of Financial Econometrics*, 15(4):649–677.
- Kosorok, M. R. (2008). *Introduction to Empirical Processes and Semiparametric Inference*. Springer.
- Kotz, S., Johnson, N. L., Balakrishnan, N., and Johnson, N. L. (2000). *Continuous Multivariate Distributions*. Wiley, New York.
- Künsch, H. R. (1989). The Jackknife and the Bootstrap for General Stationary Observations. *Annals of Statistics*, 17(3):1217–1241.
- Marquering, W. and Verbeek, M. (2004). The Economic Value of Predicting Stock Index Returns and Volatility. *Journal of Financial and Quantitative Analysis*, 39(2):407–429.
- McAleer, M. and Da Veiga, B. (2008). Single-Index and Portfolio Models for Forecasting Value-at-Risk Thresholds. *Journal of Forecasting*, 27(3):217–235.
- McCracken, M. W. (2020). Diverging tests of equal predictive accuracy. *Econometrica*, 88(4):1753–1754.
- Newey, W. K. and West, K. D. (1987). A Simple, Positive Semi-Definite, Heteroskedasticity and Autocorrelation Consistent Covariance Matrix. *Econometrica*, 55(3):703.

- Patton, A. J. (2020). Comparing Possibly Misspecified Forecasts. *Journal of Business and Economic Statistics*, 38(4):796–809.
- Post, T., Potì, V., and Karabati, S. (2018). Nonparametric Tests for Superior Predictive Ability. *Available at SSRN 3251944*.
- Riskmetrics (1996). JP Morgan Technical Document.
- Santos, A. A., Nogales, F. J., and Ruiz, E. (2012). Comparing Univariate and Multivariate Models to Forecast Portfolio Value-at-Risk. *Journal of Financial Econometrics*, 11(2):400–441.
- West, K. D. (1996). Asymptotic Inference About Predictive Ability. *Econometrica*, 64(5):1067–1084.
- White, H. (2000). A Reality Check for Data Snooping. *Econometrica*, 68(5):1097–1126.
- White, H. (2001). *Asymptotic Theory for Econometricians*. Academic Press, Cambridge, MA.
- Ziegel, J. F., Krüger, F., Jordan, A., and Fasciati, F. (2020). Robust Forecast Evaluation of Expected Shortfall. *Journal of Financial Econometrics*, 18(1):95–120.