

On the Comparison of Interval Forecasts

Ross Askanazi Francis X. Diebold
Cornerstone Research University of Pennsylvania

Frank Schorfheide Minchul Shin
University of Pennsylvania University of Illinois

January 12, 2018

Abstract: We explore interval forecast comparison when the nominal confidence level is specified, but the quantiles on which intervals are based are not specified. It turns out that the problem is difficult, and perhaps unsolvable. We first consider a situation where intervals meet the Christoffersen conditions (in particular, where they are correctly calibrated), in which case the common prescription, which we rationalize and explore, is to prefer the interval of shortest length. We then allow for mis-calibrated intervals, in which case there is a calibration-length tradeoff. We propose two natural conditions that interval forecast loss functions should meet, and we show that a variety of popular approaches to interval forecast comparison fail them. Our negative results strengthen the case for abandoning interval forecasts in favor of density forecasts: Density forecasts not only provide richer information, but also can be readily compared using known proper scoring rules like the log predictive score, whereas interval forecasts cannot.

Acknowledgments: For helpful discussion we thank Ulrich Mueller, Glenn Rudebusch, and Jonathan Wright. For research support we thank the National Science Foundation and the Real-Time Data Research Center at the Federal Reserve Bank of Philadelphia. The usual disclaimer applies.

Key words: Forecast accuracy, forecast evaluation, prediction

JEL codes: C53

Contact: fdiebold@sas.upenn.edu

1 Introduction

We consider the following situation: A researcher, R, is presented with two long sequences of univariate interval forecasts and the corresponding realizations, where each interval forecast has target (nominal) coverage $(1 - \alpha) \times 100$ percent. By “long”, we mean that we need not distinguish between consistent estimators and their probability limits, effectively working in population. Knowing nothing else, R must decide which interval forecast he prefers; that is, he must rank the two. Perhaps surprisingly, a full solution to this highly-practical and simply-stated problem remains elusive. In this paper we address it.

We emphasize – and this is the source of difficulty – that R knows only that an interval’s target coverage is $(1 - \alpha)$, not the quantiles on which it is based. A target 90 percent interval, for example, *might* be equal-tailed (target probability 5 percent in each tail), *or it might not*, in which case there are uncountably infinitely many possibilities. R sees only the two sets of intervals, each with known target coverage $(1 - \alpha)$, but with each based on unknown and potentially different quantiles.

Related foundational literature dates at least to Aitchison and Dunsmore (1968) and includes Winkler (1972) and Casella et al. (1993). Additional literature includes, among others, Granger et al. (1989) who note the tradeoff between an interval’s expected length and its coverage, Giacomini and Komunjer (2005) who explore the relative assessment of conditional quantile estimates, and Gneiting and Raftery (2007) who propose proper scoring rules for interval predictions based on identical pairs of quantiles; see also Schervish (1996).

We proceed as follows. In Section 2 we treat the case where both interval forecasts meet the optimality conditions of Christoffersen (1998) (“the Christoffersen Conditions”, CC); that is, when the 0-1 hit sequences associated with the interval forecast sequences are iid Bernoulli with the correct hit probability. In that case the prescriptive optimum, which we characterize and explore in detail, is to prefer the interval of shortest length. In Section 3 we move to the richer and largely-unexplored case where one or both interval forecasts fail the CC due to mis-calibration. We consider a natural pair of conditions that interval forecast loss functions should meet, and we show that a variety of popular approaches to interval forecast comparison fail them. We conclude in Section 4.

2 The Framework, a Prescriptive Approach, and the Christoffersen Conditions (CC)

2.1 The Basic Framework

Consider a univariate time series $\{y_t\}$ and let $y_{t+h|t}$ be an h -step-ahead point forecast of y_t . The corresponding h -step-ahead forecast error is $e_{t+h|t} = y_{t+h} - y_{t+h|t}$. Similarly, denote an h -step-ahead interval forecast (“prediction interval”) by $d_{t+h|t}(1-\alpha) = [d_{t+h|t}^l(1-\alpha), d_{t+h|t}^u(1-\alpha)]$, where the $1-\alpha$ denotes target $(1-\alpha) \times 100$ percent coverage, $0 \leq \alpha \leq 1$. Let $|d_{t+h|t}(1-\alpha)|$ be the length of the interval. The corresponding h -step “hit sequence” is $1\{y_{t+h} \in d_{t+h|t}\}$, where $1\{\cdot\}$ is the indicator function. Finally, denote an h -step-ahead density forecast by $f_{t+h}(y|\Omega_t)$, where Ω_t is the information set on which the forecast is based. The corresponding probability integral transform (PIT) is $z_{t+h|t} = F_{t+h}(y_{t+h}|\Omega_t)$, where $F_{t+h}(\cdot|\Omega_t)$ is the c.d.f. associated with the density $f_{t+h}(\cdot|\Omega_t)$. We assume that the researcher R is presented with two (or more) forecasts, which could be distinguished in terms of notation by adding j superscripts to the previously-defined objects $y_{t+h|t}$, $e_{t+h|t}$, $d_{t+h|t}$, $F_{t+h}(\cdot)$, $z_{t+h|t}$, and Ω_t .

Necessary conditions for forecast optimality with respect to an information set exist and are widely-used for all of point, interval, and density forecasts. The basic idea is just the well-known “orthogonality condition:” the relevant notion of “forecast error” should be unforecastable using information available when the forecast was made. In the canonical one-step-ahead case this is commonly assessed as: (a) point forecast errors $e_{t+1|t}$ are independently and identically distributed (iid) with mean 0; (b) interval forecast hits $1\{y_{t+h} \in d_{t+h|t}\}$ are iid with mean $1-\alpha$ as in Christoffersen (1998); and (c) density forecast PIT’s $z_{t+1|t}$ are iid $U(0,1)$ as in Diebold et al. (1998). In the general h -step-ahead case we simply replace “iid” with “at most $(h-1)$ -dependent”.

However, quite apart from “absolute evaluation,” that is, evaluation of whether a forecast satisfies conditions necessary for efficiency with respect to an information set, there is also the issue of “relative evaluation” (comparison) – simply quantifying a forecast’s performance and comparing it to that of competitors, regardless of whether any forecast under consideration is “optimal” in any sense. Two forecasts, for example, may each be efficient (optimal), but with respect to different information sets, or both may simply be inefficient. Relative evaluation is typically done for point forecasts by invoking a loss function (e.g., quadratic loss) and examining the corresponding realized average loss (e.g., mean-squared error (MSE)). In precise parallel, for density forecasts one typically invokes a loss function, e.g., the log

predictive score (LPS)

$$LPS = \log (f_{t+h}(y|\Omega_t)),$$

which is the log of the predictive density evaluated at the corresponding realization, and one examines the corresponding realized average LPS, as in Amisano and Geweke (2018) and many of the references therein. But it is much less clear what to do – and there is no “typical” procedure – for interval forecasts, to which we now turn.

2.2 The Prescriptive Optimal Interval Forecast

We begin by providing a decision-theoretic framework to derive an optimal interval forecast.¹ Denote an interval forecast by $d = [d^l, d^u]$, with length $|d| = d^u - d^l$. To reduce notational clutter, we drop the “ $t + h|t$ ” and “ $(1 - \alpha)$ ” notation, as meaning will be clear from the context. Consider a game between a forecaster F and an adversary A, where F chooses d and A chooses a scalar $\delta \in [-\infty, 0]$. Suppose that F’s loss function is

$$\mathcal{L}_F(y, d, \delta; \alpha) = |d| + \delta(1\{y \in d\} - (1 - \alpha)). \quad (1)$$

This just says that, other things the same, F prefers short intervals (small $|d|$ promotes small \mathcal{L}) and intervals that contain the realization ($y \in d$ implies that $[1\{y \in d\} - (1 - \alpha)]$ is positive, which promotes small \mathcal{L}). Note that F’s loss function reflects a tradeoff: the wider the interval, the more likely it is to contain the truth, which decreases loss as per the second term in F’s loss function, but increases loss as per the first term in the loss function.

Suppose also that A’s loss function is

$$\mathcal{L}_A(y, d, \delta; \alpha) = -\delta(1\{y \in d\} - (1 - \alpha)). \quad (2)$$

A’s loss function is very simple: he incurs negative loss (that is, positive utility) of $\delta(1 - \alpha)$ when F “gets it wrong” (i.e., when the realization is outside the interval), and positive loss of $-\delta\alpha$ when F “gets it right” (the realization is inside the interval). Hence A simply wants F to “get it wrong”, which makes explicit the sense in which A is an adversary.

Notice that d and δ appear in the loss functions of both F and A. But F chooses d and treats δ as fixed (it is set externally by A), whereas A chooses δ and treats d as fixed (it is set externally by F). Finally, assume that both players know both loss functions and the associated strategic considerations, and that they use the same posterior to calculate

¹The subsequent exposition builds on Herbst and Schorfheide (2015).

posterior expected loss (risk). The equilibrium of the game is easily obtained by considering F’s three choices:

1. Choose a correctly-calibrated interval, i.e., an interval such that $\mathbb{P}(y \in d) = (1 - \alpha)$, where \mathbb{P} denotes posterior probability. Risk to F associated with the second term of his loss function then vanishes, regardless of the adversary’s action, so his risk is governed entirely by posterior expected interval length.
2. Choose an interval such that $\mathbb{P}(y \in d) > (1 - \alpha)$. That improves risk to F associated with the second term of his loss function if $\delta < 0$, so the adversary will choose $\delta = 0$, which F knows. Hence F has no incentive to choose, and will not choose, intervals with $\mathbb{P}(y \in d) > (1 - \alpha)$.
3. Choose an interval such that $\mathbb{P}(y \in d) < (1 - \alpha)$. That worsens risk to F associated with the second term of his loss function if $\delta < 0$, so the adversary will choose $\delta = -\infty$, which F knows. Hence F has no incentive to choose, and will not choose, intervals with $\mathbb{P}(y \in d) < (1 - \alpha)$.

Clearly, then, F always prefers the first option, a correctly-calibrated interval, in which case his expected loss (risk) collapses to²

$$\mathbb{E}[\mathcal{L}_F^*(y, d; \alpha)] = \mathbb{E}[|d|], \tag{3}$$

so he prefers the interval with the shortest length. Hence the decision-theoretic prescription optimizes the length-coverage tradeoff in a very special, lexicographic, way: restrict attention to correctly-calibrated intervals, and then pick the (on average) shortest.

2.3 Characterization of the “Shortest Well-Calibrated Interval”

It will prove useful to characterize the shortest well-calibrated (sc) interval with respect to the predictive distribution of y . This interval minimizes the expected length subject to a coverage rate constraint:

$$\min_{d^l, d^u} (d^u - d^l) \quad \text{s.t.} \quad \mathbb{P}(y \in d) = 1 - \alpha. \tag{4}$$

²The * superscript indicates that the loss is conditional on the best response of the adversary A.

Because we restrict our attention to connected intervals, the constraint can be written as

$$\mathbb{P}(y < d^u) - \mathbb{P}(y < d^l) = 1 - \alpha. \quad (5)$$

The optimal interval satisfies the condition

$$f(d_{sc}^l(1 - \alpha)) = f(d_{sc}^u(1 - \alpha)), \quad (6)$$

where $f(\cdot)$ is the predictive density. (6) is obtained from the first-order conditions of the constrained optimization problem. Thus, the optimal interval equalizes the heights of the predictive density $f(y)$ at d_{sc}^l and d_{sc}^u . If the predictive density is unimodal, then the solution to the optimization problem yields the highest-probability-density interval – meaning that it includes all values y such that $f(y) \geq f(d_{sc}^l(1 - \alpha))$ – which is prominently featured in Bayesian analysis; see, for instance, the textbook by Robert (1994).

2.4 Interval Forecast Comparison Under the CC

Provided that two interval forecasts satisfy the CC either “exactly” with a hit probability equal to $1 - \alpha$ or “weakly” with a hit probability that is no less than $1 - \alpha$, the conventional prescription is to rank the intervals based on their average length, preferring shorter to longer, as intervals shorter in expectation presumably condition on more valuable information sets. This prescription is a consequence of the decision-theoretic framework presented in Section 2.2. We could simply replace the forecaster “F” by the researcher “R” to formally justify this prescription.

Unfortunately, adopting the framework of Section 2.2 has some rather undesirable implications and does not solve the interval-forecast-comparison problem in a satisfactory manner. First, forecasters receive no credit for reporting a somewhat longer interval that exceeds the nominal coverage probability. Second, and more importantly, the penalty for undercoverage is unreasonably harsh. The risk score associated with an interval that does not satisfy the coverage probability constraint is $\mathbb{E}[\mathcal{L}_R^*(y, d; \alpha)] = -\infty$. Thus, any forecast that does not achieve the nominal coverage should be disregarded and the approach produces no useful ranking of two interval forecasts whose actual coverage probabilities fall short of the nominal coverage probability, even though in practice one would prefer a forecast that yields an average length of one and a coverage probability of 92% to a forecast that has an average length of two and a coverage probability of 90% even if the nominal coverage probability is

95%. Thus, in the remainder of this paper we will explore alternative evaluation approaches.

3 Comparing Potentially Mis-Calibrated Interval Forecasts

A major practical issue is what to do when we acknowledge that, in general, the CC *fail*.³ The “shortest well-calibrated interval” result provides a prescriptive recommendation for *constructing* interval forecasts, not a descriptive recommendation for *evaluating* interval forecasts, precisely because in general the CC fail (because one or both intervals may be mis-calibrated), in which case there is a calibration-length tradeoff. In this section we study aspects of that tradeoff.⁴

3.1 Loss-Function Desiderata

We consider loss functions of the form $\mathcal{L}(d, y; \lambda)$, where the first two arguments are the endpoints of the interval forecast, $d = [d^l, d^u]$, the third argument is a realization of Y whose density function is $f(y)$ and distribution function is $F(y)$, and λ is a finite-dimensional parameter vector which includes α .⁵

There are two key desirable traits of such interval-forecast loss functions in our environment, in which the researcher R knows nothing about the intervals apart from target coverage.

D1. \mathcal{L} should (a) reflect an inverse tradeoff between length and coverage; (b) avoid the Casella paradox; and (c) have corresponding risk minimized by the shortest well-calibrated interval.

Remarks: (a) Both short expected length and correct coverage are desirable. Hence a loss function’s “indifference curves” between average length and (absolute) coverage distortion should be negatively sloped. One should be willing, for example, to accept an increase in expected length in exchange for a “large enough” improvement in coverage.

³We assume the *iid* part of the CC throughout; hence when we speak of failure of the CC we mean $\mathbb{E}[1\{y_{t+h} \in d_{t+h|t}\}] \neq 1 - \alpha$.

⁴We assume throughout that all intervals are connected, and closely related, that the predictive distribution of y , $f(y)$, has a single mode.

⁵Alternatively, we could of course cast the discussion in terms of predictive score functions. The score function $\mathcal{S}(d, y; \lambda)$ is simply $-\mathcal{L}(d, y; \lambda)$.

(b) Loss functions should, however, avoid the Casella paradox. The Casella paradox refers to situations where the obviously intuitive optimal interval is dominated in expected loss by arbitrarily short mis-calibrated intervals. It emerges when loss functions place “too much” weight on length as opposed to coverage.⁶

(c) Imagine that one of the forecaster’s is endowed with knowledge of the DGP and hence the “true” conditional distribution, denoted by $f_0(y_{t+h}|\mathcal{F}_t)$. Moreover, suppose that this forecaster reports the shortest well-calibrated interval described in Section 2.3, then, given the nominal coverage probability $1 - \alpha$, and despite the length-coverage trade-off encoded in the loss function, ideally no other forecast should yield a lower risk. ■

D2. \mathcal{L} should be evaluable with no knowledge of (a) the data generating process (DGP), or (b) the specific quantiles used by the forecasters to form the intervals.

Remark: At some level, D2 is self-evident. *Of course* evaluation of \mathcal{L} should not require knowledge of the DGP (no one knows the DGP), and *of course* it should not require knowledge of the specific left and right quantiles used (the forecasters are simply asked to provide $(1 - \alpha)$ % intervals). More formally, what we mean is that the parameter vector of the loss function, λ , should not depend on the probability measure \mathbb{P} associated with the DGP, which is used to compute expected losses $\mathbb{E}[\mathcal{L}(d, y; \lambda)]$. ■

In what follows we will restrict attention to prominent loss functions that satisfy D1(a) and D1(b). However, we will show that in general it is difficult to reconcile D1(c) and D2.

3.2 Candidate Approaches

In view of the desiderata, we now examine several interval-forecast evaluation approaches. We begin with a modification of the framework in Section 2.2 and then proceed to alternative loss/risk functions that have been proposed in the literature. A key problem for many of the loss functions is the following. Had the forecasters known that they will be evaluated based on such a loss function, they would not have reported the shortest $1 - \alpha$ coverage interval under their subjective beliefs about y_{t+h} conditional on their Ω_t information sets. Thus, using a terminology from the density forecast evaluation literature, the evaluation procedures based on such loss functions are not proper.

⁶In particular, loss functions that embed a linear length-coverage tradeoff typically fall victim to the Casella paradox. We provide an example and additional discussion in Appendix A.

3.2.1 The Two-Player Game – Revisited

The decision-theoretic framework outlined in Section 2.2 arguably fails desideratum D1(a). It does not generate a (reasonable) trade-off between average length and coverage probability. However, a small modification of the setup can generate some improvements. Suppose we restrict the choice set of the adversary A to $\delta \in [-\underline{\delta}, 0]$. In the context of this framework, we could define the vector $\lambda = [\underline{\delta}, \alpha]'$, which does not depend on the DGP and therefore satisfies D2. As under the original setup, a forecaster does not receive any credit for providing an interval with a coverage probability that exceeds $1 - \alpha$. However, there is now a trade-off between length and coverage for intervals with an actual coverage probability that falls short of the nominal coverage probability:

$$\mathbb{E}[\mathcal{L}_R^*(y, d; \lambda)] = \mathbb{E}[|d|] + \begin{cases} 0 & \text{if } \mathbb{E}[1\{y \in d\}] \geq 1 - \alpha \\ \underline{\delta}((1 - \alpha) - \mathbb{E}[1\{y \in d\}]) & \text{otherwise.} \end{cases} \quad (7)$$

As long as $\underline{\delta}$ is not too small, the Casella paradox can be avoided. Moreover, sufficiently large values of $\underline{\delta}$ may not distort the forecasters' incentive to truthfully reveal $1 - \alpha$ coverage sets based on their predictive distributions. A practical difficulty is the choice of the parameter $\underline{\delta}$. To facilitate this choice in an *ex post* evaluation, it may be convenient to replace $|d|$ by $\log |d|$ and focus on percentage differences between average interval lengths across forecasts.⁷ In this case, setting $\underline{\delta} = 1$ implies the following: compared to the shortest interval with nominal coverage probability, an interval that has the same length and undercovers by one percentage point receives the same penalty as an interval that attains the nominal coverage probability but is one percent longer than the baseline interval.

A key question is whether one can choose $\underline{\delta}$ sufficiently small to keep the tradeoff between expected length and coverage probability interesting and simultaneously satisfy desiderata D1(c) and D2. Knowledge about some features of the DGP is required to ensure that reducing the length of the optimal $1 - \alpha$ interval under the DGP does not lead to a reduction in the expected loss.

Example. Suppose that under the DGP the predictive distribution is $N(\mu_{t+h|t}, \sigma_{t+h|t}^2)$. Then the log length of the predictive interval is $|d| = \ln 2\sigma_{t+h|t} + \ln(\Phi_N^{-1}(1 - \alpha/2))$, where $\Phi_N^{-1}(\cdot)$ is the inverse cumulative density function of a standard Normal random variable. Suppose that $\alpha = 0.05$. Increasing the critical value from 0.05 to 0.06 (by 1 percentage point) reduces the log length of the interval by $\ln(1.8808/1.96) \approx -0.04$, meaning the length shrinks by

⁷Of course, from an *ex ante* perspective, this would create an incentive to report intervals that have length zero.

about 4%. This means that at $\alpha = 0.05$ if we set $\underline{\delta} \geq 5$, say, then intervals that undercover yield a larger loss. ■

3.2.2 Winkler Loss and its Relatives

The most commonly-used interval-forecast loss function is due to Winkler (1972); see also Schervish (1996):

$$\mathcal{L}_{eq}(y, d; \lambda) = |d| + \lambda_l(d^l - y)1\{y < d^l\} + \lambda_u(y - d^u)1\{y > d^u\}, \quad (8)$$

where $\lambda = [\lambda_l, \lambda_u]'$ with $\lambda_l > 0$ and $\lambda_u > 0$. The first term penalizes the length of the interval, and the second and third terms penalize misses from the left and right tails, respectively. If the interval misses the realization from the left ($y < d^l$), the second term ($\lambda_l(d^l - y)$) becomes positive and the third term becomes zero. If the interval misses from the right ($y > d^u$), the second term becomes zero and the third term ($\lambda_u(y - d^u)$) becomes positive.

It can be shown (Gneiting and Raftery, 2007) that the optimal interval forecast under the Winkler loss has left and right endpoints given by the $1/\lambda_l$ and $1 - 1/\lambda_u$ quantiles, so that for the evaluation of intervals with $(1 - \alpha)$ coverage, one should set $1/\lambda_l + 1/\lambda_u = \alpha$. A popular choice for λ_l and λ_u in practice is $\lambda_l = \lambda_u = 2/\alpha$, in which case the optimal interval is equal-tailed (Gneiting and Raftery, 2007),

$$\mathcal{L}(d, y; \alpha) = |d| + \frac{2}{\alpha} \inf_{x \in d} |y - x|.$$

Notice, crucially, that *good interval forecasting under Winkler loss amounts to good quantile forecasting, for known, stated, quantiles.*

Some other well-known loss functions are very similar to Winkler's. For example, Schmalensee (1976) proposes

$$\mathcal{L}_S(y, d, \lambda_l, \lambda_u) = |d| + \lambda_l(d^l - y)1\{y < d^l\} + \lambda_u(y - d^u)1\{y > d^u\} + \left| y - \frac{d^u + d^l}{2} \right|. \quad (9)$$

The first three terms are the same as Winkler's. The difference is an additional penalty for the divergence between the realization and the midpoint of the interval. Nevertheless, good interval forecasting under this loss function amounts to good quantile forecasting.

The key insight is that all Winkler-type loss functions are quantile-type loss functions – they involve properties of estimates of known quantiles that define the left and right interval endpoints; see also Gneiting (2011b). Quantile-type loss functions are fine if one knows (or is

willing to assume) that the intervals to be compared are equal-tailed, or, more generally, are based on known left and right quantiles. But that’s not our case – we are simply provided with two sequences of intervals, and we must decide which we prefer.

The quantile-type loss functions are unable to resolve the tension between desiderata D1(c) and D2. Consider Winkler loss. In order to generate a $1 - \alpha$ interval forecast, it has to be the case that $1/\lambda_l + 1/\lambda_u = \alpha$. Moreover, in order to satisfy D2, the parameters λ_l and λ_u cannot depend on features of the DGP. If it were known that the DGP is of the form $y_{t+h}|\mathcal{F}_t \sim N(\mu_{t+h|t}, \sigma_{t+h|t}^2)$, then, in every period t , the optimal $1 - \alpha$ interval forecast would correspond to the $\alpha/2$ and $1 - \alpha/2$ quantiles, which would also minimize the expected Winkler loss for $\lambda_l = \lambda_u = 2/\alpha$. Thus, no further information about the DGP is required to satisfy D1(c) and D2. However, if the DGP is a mixture of normal distributions with time-varying mixture weights, then it is generally not true that one can find λ ’s such that the optimal $1 - \alpha$ interval forecast also minimizes the Winkler loss. Thus, D1(c) is violated.

3.2.3 Length / Hit Loss

Loss functions that directly respect a length / hit tradeoff are of the form,

$$\mathcal{L}_G(d, y; \lambda) = M(|d|, 1\{y \in d\}; \lambda), \quad (10)$$

where $M(\cdot, \cdot)$ is a potentially nonlinear function parametrized by $\lambda \in \Lambda \subseteq \mathcal{R}^n$. The first argument is the interval length and the second argument is the hit indicator.⁸ $M(\cdot, \cdot)$ is increasing in length, and $M(x, 1; \lambda) < M(x, 0; \lambda)$ for all x and $\lambda \in \Lambda$. Note that y enters M only through the hit indicator. We assume that M ’s risk minimizer is unique, and that $M(\cdot, 1; \lambda)$ and $M(\cdot, 0; \lambda)$ are smooth functions (twice differentiable).⁹

One prominent example is the following loss function proposed by Casella et al. (1993),

$$\mathcal{L}_{CHR}(d, y; \lambda) = \frac{|d|}{|d| + \lambda} - 1\{y \in d\} \text{ for } \lambda \in (0, \infty), \quad (11)$$

and another is the “most likely interval” loss of Schlag and van der Weele (2015),

$$\mathcal{L}_{MLI}(d, y; \lambda) = - \left(1 - \frac{|d|}{b - a}\right)^\lambda 1\{y \in d\} \text{ for } \lambda \in (0, \infty), \quad (12)$$

⁸The function $\mathcal{L}_F(y, d, \delta; \alpha)$ in Section 2.2 is a function of the length and the hit indicator, but it also depends on the decision δ of the adversary.

⁹Casella et al. (1993) provide a set of conditions for M that avoid the pathological Casella’s paradox. See Appendix A for more details.

where the support of y is restricted to a bounded interval, $[a, b]$.

The bounds of the risk-minimizing interval d_* satisfy the following for any $\lambda \in \Lambda$:

$$f(d_*^l(\lambda)) = f(d_*^u(\lambda)). \quad (13)$$

To see this, write $|d| = d^u - d^l$ and note that the expected loss is

$$\mathbb{E}[\mathcal{L}_G(d, y; \lambda)] = M(|d|, 1; \lambda)[F(d^u) - F(d^l)] + M(|d|, 0; \lambda)[1 - F(d^u) + F(d^l)], \quad (14)$$

so that the first order conditions with respect to d^l and d^u are

$$\begin{aligned} \partial d^l : & -M'(|d|, 1; \lambda)[F(d^u) - F(d^l)] - M(|d|, 1; \lambda)f(d^l) \\ & - M'(|d|, 0; \lambda)[1 - F(d^u) + F(d^l)] + M(|d|, 0; \lambda)f(d^l) = 0 \\ \partial d^u : & M'(|d|, 1; \lambda)[F(d^u) - F(d^l)] + M(|d|, 1; \lambda)f(d^l) \\ & + M'(|d|, 0; \lambda)[1 - F(d^u) + F(d^l)] - M(|d|, 0; \lambda)f(d^u) = 0. \end{aligned} \quad (15)$$

After rearranging terms, we find that $f(d^l) = f(d^u)$. This implies that the optimal interval prediction under loss function (10) satisfies one of the requirements for the shortest well-calibrated interval, as stated in equation (6).

Unfortunately, however, the coverage rate of this interval depends on λ . That is, λ needs to be set in a way that generates $(1 - \alpha)$ coverage, which requires knowledge of the true predictive distribution (i.e., the DGP, which is specified as a conditional distribution). In time-series settings, moreover, the shape of the predictive distribution can change over time, which implies that the λ that induces $(1 - \alpha)$ coverage can also vary over time. In any event, it is clear that D1(c) and D2 cannot be satisfied simultaneously.

3.2.4 Direct Risk-Based Comparisons

In practice, researchers and practitioners often separately report and examine interval forecasts' average length and empirical coverage, as in Granger et al. (1989). In this way, people can presumably rank interval forecasts according to their own preferences (i.e., relative importance of sharpness vs. calibration).

An interesting question is therefore whether there exists a loss function whose expected loss (risk) (a) can be written as a function of expected length and coverage, and (b) is

minimized by the shortest well-calibrated interval. Requirement (a) allows us to compute the expected loss easily by just replacing the expected length and coverage rate by their sample counterparts.

This motivates the direct study of risk functions, where risk depends only on expected length ($L = \mathbb{E}[|d|]$) and the difference between the target coverage rate and the actual coverage rate ($C = (1 - \alpha) - \mathbb{E}[1\{y \in d\}]$, the “coverage distortion”). The risk function (7) that resulted from the two player game takes this particular form. More generally, consider

$$\mathbb{E}[\mathcal{L}_{G'}(d, y; \lambda)] = M(L, C; \lambda), \quad (16)$$

where the expectation is taken with respect to the true predictive distribution, the function M is smooth (piecewise differentiable) for $\lambda \in \Lambda \subseteq \mathcal{R}^n$, and

$$\frac{\partial M(L, C; \lambda)}{\partial L} > 0 \quad \text{and} \quad \frac{\partial M(L, C; \lambda)}{\partial C} > 0, \quad (17)$$

whenever the derivatives exist.¹⁰

Example 1 (Linear Risk). An immediate example is

$$\mathbb{E}[\mathcal{L}_{G'}(d, y; \lambda)] = L + \lambda C, \quad (18)$$

where $\lambda \in (0, \infty)$. Due to the linearity, this risk function corresponds to a loss function of the form (10). ■

Example 2 (CWC Risk). Not all risk functions that satisfy (16) correspond to loss functions that satisfy (10). For example, Khosravi and Nahavandi (2014) propose a “coverage width-based criterion” (CWC), defined as

$$\mathbb{E}[\mathcal{L}_{CWC}(d, y; \lambda)] = L + 1\{C > 0\}e^{\lambda C}, \quad (19)$$

where $\lambda \in (0, \infty)$. ■

Direct measures (16) have some good properties. In particular, under our earlier assumptions on $M(\cdot, \cdot)$, for any given λ the risk minimizer satisfies:

$$f(d_*^l(\lambda)) = f(d_*^u(\lambda)), \quad (20)$$

¹⁰As discussed in Section 3.2.1, in principle we would also need to impose restrictions on $M(\cdot, \cdot)$ sufficient to ensure avoidance of Casella’s paradox, but we do not pursue that here.

which follows trivially from the first order conditions,

$$\begin{aligned} \partial d^l : -\frac{\partial M(L, C)}{\partial L} + \frac{\partial M(L, C)}{\partial C} f(d^l) &= 0 \\ \partial d^u : \frac{\partial M(L, C)}{\partial L} - \frac{\partial M(L, C)}{\partial C} f(d^u) &= 0, \end{aligned} \tag{21}$$

rearrangement of which yields $f(d^l) = f(d^u)$. This implies that the optimal interval prediction under risk function (16) satisfies one of requirements for the shortest well-calibrated interval (equation (6)).

Unfortunately, however, direct risk measures (16) also have earlier-discussed bad properties that plagued other approaches. In particular, the coverage rate of the optimal predictive interval depends on λ . That is, for this class of risk function to be proper for the shortest connected interval with $(1 - \alpha)$ coverage, we need to select λ accordingly, but λ depends on the DGP $f_0(\cdot)$, which is unknown.

Example 3 (Linear Risk – Revisited). Consider again the linear loss function,

$$\mathbb{E}[\mathcal{L}_{G'}(d, y; \lambda)] = L + \lambda C, \tag{22}$$

where $\lambda \in (0, \infty)$. In this case $1/\lambda = f(d^l) = f(d^u)$. If f is the standard normal density function, then $\lambda = 1/f(q_{\alpha/2})$ uniquely produces the shortest connected interval with $(1 - \alpha)$ coverage, where $q_{\alpha/2}$ is the $\alpha/2$ quantile of the standard normal distribution. ■

3.2.5 On Medians, Means, and Modes

Note that the use of equal-tailed interval forecasts naturally parallels the use of conditional-median point forecasts. To see this, note that the equal-tailed interval with ϵ coverage probability approaches the conditional median as $\epsilon \rightarrow 0$. The associated absolute-error loss function is well studied in the literature, which facilitates finding good scoring rules for equal-tailed intervals.

As an aside we note that, interestingly, little attention has been given to interval forecasts that parallel conditional-*mean* point forecasts. We are aware of just one paper, Rice et al. (2008), which considers the loss function

$$\mathcal{L}(d, y; \lambda) = \lambda \frac{|d|}{2} + \frac{(y - |d^u + d^l|/2)^2}{|d|/2}, \tag{23}$$

which leads to the interval forecast $\left[m \pm s/\sqrt{\lambda} \right]$, where m and s are the mean and standard deviation of the predictive distribution of y . The coverage rate is implicitly defined through λ and the predictive distribution of y . As the coverage rate approaches zero, the interval approaches the mean of the predictive distribution.¹¹

Now, finally, consider finding good scoring rules for the shortest connected interval. Shortest connected interval forecasts naturally parallel conditional-mode point forecasts, because the shortest connected interval with ϵ coverage probability approaches the conditional mode as $\epsilon \rightarrow 0$.¹² The associated 0-1 loss function, however, is mathematically more clumsy than absolute-error loss, which impedes finding good scoring rules for shortest connected intervals. Gneiting (2011a), for example, discusses a scoring rule for the mode, but he does not find a scoring rule that selects the mode as the global optimum.

3.3 Additional Remarks

First, the overarching point is that, to use any of the leading loss or risk functions discussed above for evaluating interval forecasts, in general one needs to know at least some features of the DGP. This creates a tension between desiderata D1(c) and D2, which often can only be jointly satisfied for a limited class of DGP's

Second, shortest connected intervals have drawbacks even beyond the difficulty of finding proper scoring rules. One example: The shortest interval may not be invariant to nonlinear transformations $g(\cdot)$; hence, for example if $[a, b]$ is the shortest connected interval for y with support $y > 0$, then $[a^2, b^2]$ is not necessarily the shortest connected interval for y^2 . Another example: If the predictive density is multi-modal, the shortest interval is not necessarily the shortest connected interval; rather, it can be a union of disjoint intervals.

Third, various methods have been proposed to construct the sharpest predictive interval (the shortest connected interval for a class of unimodal distributions) based on parametric models (Hyndman, 1995, 1996), nonparametric models (Polonik and Yao, 2000), and semi-parametric models Wu (2012). However, quantile-based predictive intervals remain dominant, probably because proper scoring rules are readily available for quantile-based intervals. People seem to be willing to sacrifice some sharpness in exchange for availability of a proper scoring rule.

Fourth, we note, however, that under some additional assumptions it may be possible

¹¹Note that setting the target coverage rate requires knowledge of the predictive distribution of y .

¹²Here we are implicitly assuming unimodality of the predictive distribution, as the shortest connected interval is the high probability density region when the predictive interval is unimodal.

to infer the predictive distribution from the interval forecast, which would let us move from interval forecast comparisons to the well-understood world of density forecast comparisons using well-known proper scoring rules. The conditions are of course strong, but perhaps not absurdly strong. Suppose that each interval prediction is a shortest connected interval (based on predictive distribution \tilde{P} and density \tilde{f} , which do not necessarily match the true predictive distribution P and density f). Then the interval forecast $d = [d^l, d^u]$ contains the following information: (1) $\tilde{P}(y \in d) = 1 - \alpha$, and (2) $\tilde{f}(d^l) = \tilde{f}(d^u)$. If we knew \tilde{P} , then the log predictive score $\log \tilde{P}(y)$ could serve as a proper scoring rule, but we know only d , not \tilde{P} .

Suppose, however, that we know the true predictive distribution P . Then we could infer \tilde{P} by solving

$$\begin{aligned} \min_q \int P \log \left(\frac{P}{\tilde{P}} \right) d\mu & \quad (24) \\ \text{s.t. } \tilde{P}(y \in d) = 1 - \alpha \quad \text{and} \quad \tilde{f}(d^l) = \tilde{f}(d^u), & \end{aligned}$$

and then compute the log predictive score for the interval forecast by

$$\mathcal{L}_{KL}(y, d) = \log \tilde{P}(y; d). \quad (25)$$

This scoring rule is proper, so long as the true predictive distribution is unimodal and we consider a connected interval. To see this, suppose that a forecaster's predictive distribution is $\hat{P} = P$. Then the solution to the minimization problem is any \tilde{P} such that $\hat{P} = P$ almost everywhere, which yields a value of 0 for the objective function (24). Achieving a value of 0 implies that the log predictive score is maximized with probability 1.¹³

Fifth, in general, making operational interval comparisons with desirable properties requires (some) knowledge of the “true” predictive distribution, whereas it is in fact unknown. But one could perhaps approximate it flexibly using, for example, ARMA conditional mean dynamics and GARCH conditional variance dynamics, together with a nonparametrically-specified conditional density. This approximation could then be used to calibrate the loss and risk functions that are used to evaluate the interval forecasts.

¹³This scoring rule will generally not be strictly proper, however, in that many predictive distributions can have the same shortest predictive interval as the true distribution and will thus achieve 0 loss.

4 Conclusion

We have explored aspects of comparing two interval forecasts when the nominal confidence level is specified, but the quantiles on which intervals are based are not specified. The problem seems simple, and the lack of attention to it in the literature seems odd. It turns out that the problem is generally quite difficult, and perhaps unsolvable.

We first considered a situation where both intervals meet the Christoffersen conditions (in particular, where both are correctly calibrated), in which case a common prescription is to prefer the interval of shortest length, and we explored and rationalized that prescription. We then allowed for mis-calibrated intervals, which cannot happen under the Christoffersen conditions. We proposed two natural desiderata for interval forecast loss functions in that case, and we showed that a variety of popular approaches fail to meet them.

A positive suggestion that springs from our negative results is that interval forecasts might be beneficially abandoned, with focus instead placed on making and comparing complete density forecasts. This has already been happening in recent decades because density forecasts are much richer than interval forecasts in terms of the information conveyed, yet easily computed by simulation in both frequentist and Bayesian frameworks. Our results strengthen the case for abandoning interval forecasts in favor of density forecasts: density forecasts can be readily compared using known proper scoring rules like the log predictive score, whereas interval forecasts cannot.

Appendices

A Two Motivating Examples

A.1 Example 1: Casella's Paradox

Recently, Gneiting and Raftery (2007) consider certain aspects of our problem. They consider loss functions that trade off length and calibration, such as the canonical

$$\mathcal{L}(d, y; \lambda) = \lambda|d| - 1\{y \in d\},$$

where λ is a scale parameter chosen by a researcher the R. However, consider the simple case of forecasting a Gaussian $y \sim N(0, \sigma^2)$. Suppose forecaster 1 issues a degenerate interval that is just the set $d_1 = \{0\}$, while forecaster 2 issues the correct equal-tailed interval $d_2 = [-z_{\alpha/2}\sigma, z_{\alpha/2}\sigma]$, with $z_{\alpha/2}$ denoting the α critical value. Then the expected losses are:

$$\mathbb{E}[\mathcal{L}(d_1, y; \lambda)] = 0, \quad \mathbb{E}[\mathcal{L}(d_2, y; \lambda)] = 2\lambda t_{\alpha}\sigma - (1 - \alpha).$$

In order to avoid the selection of the degenerate interval forecast d_1 , R must set

$$\lambda < \frac{1 - \alpha}{2z_{\alpha/2}\sigma}.$$

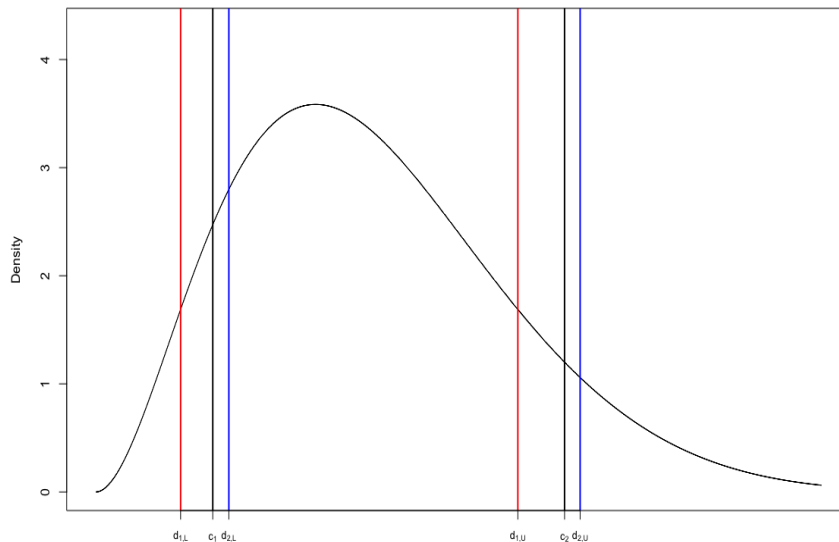
That the parameter λ depends on the possibly unknown σ , and that for a fixed λ increasing σ will inevitably yield selection of the degenerate interval instead of an increasingly wide interval, is known as ‘‘Casella’s paradox’’ (see Casella et al. (1993)). The Casella et al. (1993) solution is to augment the loss function with a size function $S(\cdot)$, so that

$$\mathcal{L}(d, y; S(\cdot)) = S(|d|) - 1\{y \in d\}$$

It is simple to derive conditions on $S(\cdot)$ that rule out the selection of *empty* intervals over those with correct coverage; however, parametric forms of $S(\cdot)$ that allow a researcher to consistently describe their desired tradeoff between coverage and length prove elusive.

Figure 1: Motivating Example

The predictive distribution is a realized beta distribution with size parameters $a = 3$, $b = 10$. Desired coverage is $1 - \alpha = .8$, i.e. $\alpha = .2$. Interval d_1 is constructed as the shortest interval with correct coverage. Interval d_2 is constructed as the equal-tailed interval, but is done so incorrectly: The lower endpoint is shifted up by some ϵ_1 , the upper endpoint is shifted up by some $\epsilon_2 \neq \epsilon_1$. c_1 is the $\alpha/2$ quantile, c_2 is the $1 - \alpha/2$ quantile. Interval d_2 will be chosen over d_1 by Gneiting and Schervish loss functions despite being longer and having incorrect coverage.



A.2 Example 2: Equal-Tailed Intervals and Asymmetric Distributions

The Gneiting and Raftery (2007) loss function

$$\mathcal{L}(d, y; \alpha) = |d| + \frac{2}{\alpha} \inf_{x \in d} |y - x|$$

is only applicable to equal-tailed intervals, and in asymmetric cases can rule out shorter intervals with correct coverage in favor of poorer intervals that are closer to the $\alpha/2$ and $1 - \alpha/2$ quantiles. Consider the following example:

- Data is generated from a Beta distribution with shape parameters $\alpha = 3$, $\beta = 10$.
- Interval d_1 is constructed as the shortest interval with correct coverage.
- Interval d_2 is constructed as the equal-tailed interval, but is done so incorrectly: The lower endpoint is shifted up by some ϵ_1 , the upper endpoint is shifted up by some

$$\epsilon_2 \neq \epsilon_1.$$

- The intervals will be ranked according to the equal-tailed interval loss function.

Visually this can be seen in Figure 1. Thus we can see that since the endpoints of d_2 are far closer to c_1, c_2 , that d_2 will be selected over d_1 even though:

- $\mathbb{P}(y \in d_1) = 1 - \alpha = .8$, while $\mathbb{P}(y \in d_2) = .78 < .8 = 1 - \alpha$.
- $|d_1| = .278$, while $|d_2| = .288$.

The crux of the issue is that the forecasters and the researcher R do not have the same target. In using this loss function, R is effectively stating that he prefers intervals closer to the given quantiles over intervals with correct coverage that are far from these quantiles. Forecaster who believe the priorities are only size and coverage will be routinely discarded, despite doing their best for their given objective.

References

- Aitchison, J. and I. Dunsmore (1968), “Linear-Loss Interval Estimation of Location and Scale Parameters,” *Biometrika*, 55, 141–148.
- Amisano, G. and J. Geweke (2018), “Prediction Using Several Macroeconomic Models,” *Review of Economics and Statistics*, 100, in press.
- Casella, G., J. Hwang, and C. Robert (1993), “A Paradox in Decision-Theoretic Interval Estimation,” *Statistica Sinica*, 3, 141–155.
- Christoffersen, P. (1998), “Evaluating Interval Forecasts,” *International Economic Review*, 39, 841–62.
- Diebold, F.X., T.A. Gunther, and A.S. Tay (1998), “Evaluating Density Forecasts, with Applications to Financial Risk Management,” *International Economic Review*, 39, 863–883.
- Giacomini, R. and I. Komunjer (2005), “Evaluation and Combination of Conditional Quantile Forecasts,” *Journal of Business and Economic Statistics*, 23, 416–431.
- Gneiting, T. (2011a), “Making and Evaluating Point Forecasts,” *Journal of American Statistical Association*, 106, 746–762.
- Gneiting, T. (2011b), “Quantiles as Optimal Point Forecasts,” *International Journal of Forecasting*, 27, 197–207.
- Gneiting, T. and A.E. Raftery (2007), “Strictly Proper Scoring Rules, Prediction, and Estimation,” *Journal of the American Statistical Association*, 102, 359–378.
- Granger, C.W.J., M. Kamstra, and H. White (1989), “Interval Forecasting: An Analysis Based Upon ARCH-quantile Estimators,” *Journal of Econometrics*, 40, 87–96.
- Herbst, E. and F. Schorfheide (2015), *Bayesian Estimation of DSGE Models*, Princeton University Press.
- Hyndman, R.J. (1995), “Highest-density Forecast Regions for Non-linear and Non-normal Time Series Models,” *Journal of Forecasting*, 14, 431–441.
- Hyndman, R.J. (1996), “Computing and Graphing Highest Density Regions,” *The American Statistician*, 50, 120–126.

- Khosravi, A. and S. Nahavandi (2014), “Closure to the Discussion of “Prediction Intervals for Short-Term Wind Farm Generation Forecasts” and “Combined Nonparametric Prediction Intervals for Wind Power Generation” and the Discussion of “Combined Nonparametric Prediction Intervals for Wind Power Generation”,” *IEEE Transactions on Sustainable Energy*, 5, 1022–1023.
- Polonik, W. and Q. Yao (2000), “Conditional Minimum Volume Predictive Regions for Stochastic Processes,” *Journal of the American Statistical Association*, 95, 509–519.
- Rice, K.M., T. Lumley, and A.A. Szpiro (2008), “Trading Bias for Precision: Decision Theory for Intervals and Sets,” Biostatistics Working Paper 336, University of Washington.
- Robert, Christian P. (1994), *The Bayesian Choice*, Springer Verlag.
- Schervish, M.J. (1996), *Theory of Statistics*, Springer.
- Schlag, K.H. and J.J. van der Weele (2015), “A Method to Elicit Beliefs As Most Likely Intervals,” *Judgment and Decision Making*, 10, 456–468.
- Schmalensee, R. (1976), “An Experimental Study of Expectation Formation,” *Econometrica*, 44, 17–41.
- Winkler, R.L. (1972), “A Decision-Theoretic Approach to Interval Estimation,” *Journal of the American Statistical Association*, 67, 187–191.
- Wu, J.J. (2012), “Semiparametric Forecast Intervals,” *Journal of Forecasting*, 31, 189–228.