

Comparing Possibly Misspecified Forecasts*

Andrew J. Patton
Duke University

This version: 24 August 2018.

Abstract

Recent work has emphasized the importance of evaluating estimates of a statistical functional (such as a conditional mean, quantile, or distribution) using a loss function that is consistent for the functional of interest, of which there are an infinite number. If forecasters all use correctly specified models free from estimation error, and if the information sets of competing forecasters are nested, then the ranking induced by a single consistent loss function is sufficient for the ranking by any consistent loss function. This paper shows, via analytical results and realistic simulation-based analyses, that the presence of misspecified models, parameter estimation error, or nonnested information sets, leads generally to sensitivity to the choice of (consistent) loss function. Thus, rather than merely specifying the target functional, which narrows the set of relevant loss functions only to the *class* of loss functions consistent for that functional, forecast consumers or survey designers should specify the single specific loss function that will be used to evaluate forecasts. An application to survey forecasts of US inflation illustrates the results.

Keywords: Survey forecasts, economic forecasting, point forecasting, model misspecification, Bregman distance, proper scoring rules, consistent loss functions.

J.E.L. codes: C53, C52, E37.

AMS 2010 Classifications: 62M20, 62P20.

*For helpful comments I thank Tim Bollerslev, Dean Croushore, Frank Diebold, Tilmann Gneiting, Jia Li, Robert Lieli, Minchul Shin, Allan Timmermann and seminar participants at Boston College, Columbia, Duke, Penn, Princeton, St. Louis Federal Reserve, 8th French Economics conference, NBER Summer Institute, Nordic Econometric Society meetings, and the World Congress of the Econometric Society. Contact address: Department of Economics, Duke University, 213 Social Sciences Building, Box 90097, Durham NC 27708-0097. Email: andrew.patton@duke.edu.

1 Introduction

Misspecified models pervade the observational sciences and social sciences. In such fields, researchers must contend with limited data, which inhibits both their ability to refine their models, thereby introducing the risk of model misspecification, and their ability to estimate these models with precision, introducing estimation error (parametric or nonparametric). This paper considers the implications of these empirical realities for the comparison of forecasts, in light of recent work in statistical decision theory on the importance of the use of consistent scoring rules or loss functions in forecast evaluation, see Gneiting (2011a). This paper shows that in analyses where forecasts are possibly based on models that are misspecified, subject to estimation error, or that use nonnested information sets (e.g., expert forecasters using different proprietary data sets), the choice of scoring rule or loss function is even more critical than previously noted.

Recent work in the theory of prediction has emphasized the importance of the choice of loss function used to evaluate the performance of a forecaster. In particular, there is a growing recognition that the loss function used must “match,” in a specific sense clarified below, the quantity that the forecaster was asked to predict, for example the mean, the median, or the probability of a particular outcome (e.g., rain, a recession), etc. In the widely-cited “Survey of Professional Forecasters,” conducted by the Federal Reserve Bank of Philadelphia, experts are asked to predict a variety of economic variables, with questions such as “What do you expect to be the annual average CPI inflation rate over the next 5 years?” In the Thomson Reuters/University of Michigan Survey of Consumers, respondents are asked “By about what percent do you expect prices to go (up/down) on the average, during the next 12 months?” The presence of the word “expect” in these questions is an indication (at least to statisticians) that the respondents are being asked for their mathematical expectation of future inflation. The oldest continuous survey of economists’ expectations, the Livingston survey, on the other hand, simply asks “What is your forecast of the average annual rate of change in the CPI?” leaving the specific type of forecast unstated.

In point forecasting, a loss function is said to be “consistent” for a given statistical functional (e.g., the mean, median, etc.), if the expected loss is minimized when the given functional is used as

the forecast, see Gneiting (2011a) and discussion therein. For example, a loss function is consistent for the mean if no other quantity leads to a lower expected loss than the mean. The class of loss functions that is consistent for the mean is known as the Bregman class, see Savage (1971), Banerjee *et al.* (2005) and Bregman (1967), and includes the squared-error loss function as a special case. The class of loss functions that is consistent for the α -quantile is known as the generalized piecewise linear (GPL) class, see Gneiting (2011b), which nests the familiar piece-wise linear function from quantile regression, see Koenker *et al.* (2017) for example. In density or distribution forecasting the analogous idea is that of a “proper” scoring rule, see Gneiting and Raftery (2007): a scoring rule is proper if the expected loss under distribution P is minimized when using P as the distribution forecast. Evaluating estimates of a given functional using consistent loss functions or proper scoring rules is a minimal requirement for sensible rankings of the competing forecasts.

Gneiting (2011a, p757) summarizes the implications of the above work as follows: *“If point forecasts are to be issued and evaluated, it is essential that either the scoring function be specified ex ante, or an elicitable target functional be named, such as the mean or a quantile of the predictive distribution, and scoring functions be used that are consistent for the target functional.”* This paper contributes to the literature by refining this recommendation to reflect real-world deviations from the ideal predictive environment, and suggests that only the first part of the above recommendation should stand; specifying the target functional is generally *not* sufficient to elicit a forecaster’s best (according to a given, consistent, loss function) prediction. Instead, forecasters should be told the single, specific loss function that will be used to evaluate their forecasts.

Firstly, I show that when two competing forecasts are generated using models that are correctly specified, free from estimation error, and when the information sets of one of the forecasters nests the other, the ranking of these forecasts based on a single consistent loss function is sufficient for their ranking using *any* consistent loss function (subject of course to integrability conditions). This is established for the problem of mean forecasting, quantile forecasting (nesting the median as a special case), and distribution forecasting.

Secondly, and with more practical importance, I show via analytical and realistic numerical examples that when *any* of these three conditions is violated, i.e. when the competing forecasts are

based on nonnested information sets, or misspecified models, or models with estimated parameters, the ranking of the forecasts is generally sensitive to the choice of consistent loss function. This result has important implications for survey forecast design and for forecast evaluation more generally.

I illustrate the ideas in this paper with a study of the inflation forecasting performance of respondents to the Survey of Professional Forecasters (SPF) and the Michigan Survey of Consumers, as well as the Federal Reserve staff’s “Greenbook” forecasts. Under squared-error loss, I find that the Greenbook forecast beats SPF, which in turn beats Michigan, but when a Bregman loss function is used that penalizes over- or under-predictions more heavily, the rankings of these forecasts switches. I also consider comparisons of individual respondents to the SPF, and find cases where the ranking of two forecasters is sensitive to the particular choice of Bregman loss function, and cases where the ranking is robust across a range of Bregman loss functions.

The (in)sensitivity of rankings to the choice of loss function also has implications for the use of multiple loss functions to compare a given collection of forecasts. If the loss functions used are not consistent for the same statistical functional, then it is not surprising that the rankings may differ across loss functions, see Engelberg *et al.* (2009), Gneiting (2011a) and Patton (2011). If the loss functions are consistent for the same functional, then in the absence of misspecified models, estimation error or nonnested information sets, the results in this paper show that using multiple measures of accuracy adds no information beyond using just one measure. (Note, however, that loss functions may have different sampling properties, and a judicious choice of loss function may lead to improved efficiency.) In the presence of these real-world forecasting complications, averaging the performance across multiple measures could mask true out-performance under one specific loss function. In recent work, Ehm *et al.* (2016) obtain mixture representations for the classes of loss functions consistent for quantiles and expectiles which can be used to determine whether one forecast outperforms another across all consistent loss functions.

This paper is related to several recent papers. Elliott *et al.* (2016) study the problem of forecasting binary variables with binary forecasts, and the evaluation and estimation of models based on consistent loss functions. Merkle and Steyvers (2013) also consider forecasting binary variables, and provide an example where the ranking of forecasts is sensitive to the choice of

consistent loss function. Lieli and Stinchcombe (2013, 2017) study the identifiability of a forecaster’s loss function given a sequence of observed forecasts, and find in particular for discrete random variables that whether the forecast is constrained to have the same support as the target variable or not has crucial implications for identification. In particular, Bregman losses become identifiable (up to scale) under such restrictions, while GPL losses are still observationally equivalent. Holzmann and Eulert (2014) show that (correctly-specified) forecasts based on larger information sets lead to lower expected loss. I build on these works, and the important work of Gneiting (2011a), to show the strong conditions under which the comparison of forecasts is insensitive to the choice of loss function. A primary goal of this paper is to show that in many realistic prediction environments, sensitivity to the choice of consistent loss function is the norm, not the exception.

A concrete outcome of this paper is the following. In macroeconomic forecasting, mean squared error (MSE) and mean absolute error (MAE) are popular ways to compare forecast accuracy, see Elliott and Timmermann (2016) for example. If the target variable is known to be symmetrically distributed, then the rankings by MSE and MAE will be the same, in the limit, if the forecasts being compared are based on nested information sets, and are free from both estimation error and model misspecification. However, if any of these ideal conditions are violated then the rankings yielded of MSE and MAE need not be the same, and the choice of loss function will affect the ranking. Similarly, in volatility forecasting MSE and QLIKE (see equation 5 below) are widely used in forecast comparisons, e.g. see Bauwens *et al.* (2012). These are both members of the Bregman family of loss functions, and so in the ideal forecasting environment they will yield, asymptotically, the same rankings of volatility forecasts. However, outside of the ideal environment rankings will generally be sensitive to the choice of loss function.

The remainder of the paper is structured as follows. Section 2 presents positive and negative results on forecast comparison in the absence and presence of real-world complications like nonnested information sets and misspecified models, covering mean, quantile and distribution forecasts. Section 3 considers realistic simulation designs that illustrate the main ideas of the paper, and Section 4 presents an analysis of US inflation forecasts. The appendix presents the main proofs, and a supplemental web appendix contains additional proofs and derivations.

2 Comparing forecasts using consistent loss functions

2.1 Mean forecasts and Bregman loss functions

The most well-known loss function is the quadratic or squared-error loss function:

$$L(y, \hat{y}) = (y - \hat{y})^2 \tag{1}$$

Under quadratic loss, and given standard regularity conditions, the optimal forecast of a variable Y_t is well-known to be the (conditional) mean:

$$\hat{Y}_t^* \equiv \arg \min_{\hat{y} \in \mathcal{Y}} \mathbb{E}[L(Y_t, \hat{y}) | \mathcal{F}_t] \tag{2}$$

$$= \mathbb{E}[Y_t | \mathcal{F}_t], \text{ if } L(y, \hat{y}) = (y - \hat{y})^2 \tag{3}$$

where \mathcal{F}_t is the information set available to the forecaster for predicting Y_t , and \mathcal{Y} is the set of possible forecasts of Y_t , which is assumed to be at least as large as the support of Y_t . More generally, the conditional mean is the optimal forecast under any loss function belonging to a general class of loss functions known as Bregman loss functions (see Banerjee *et al.*, 2005 and Gneiting, 2011a). The class of Bregman loss functions is then said to be “consistent” for the (conditional) mean functional. Elements of the Bregman class of loss functions, denoted $\mathcal{L}_{Bregman}$, take the form:

$$L(y, \hat{y}) = \phi(y) - \phi(\hat{y}) - \phi'(\hat{y})(y - \hat{y}) \tag{4}$$

where $\phi : \mathcal{Y} \rightarrow \mathbb{R}$ is any strictly convex function. (Here and throughout, we will focus on *strict* consistency of a loss function, which in this section requires strict convexity of ϕ ; see Gneiting (2011) for discussion of consistency versus strict consistency.) Moreover, this class of loss functions is also *necessary* for conditional mean forecasts, in the sense that if the optimal forecast is known to be the conditional mean, then it must be that the forecast was generated by minimizing the expected loss of some Bregman loss function. Two prominent examples of Bregman loss functions are quadratic loss (equation (1)) and QLIKE loss (Patton, 2011), which is applicable for strictly positive random variables:

$$L(y, \hat{y}) = \frac{y}{\hat{y}} - \log \frac{y}{\hat{y}} - 1 \tag{5}$$

The quadratic and QLIKE loss functions are unique (up to location and scale constants) in that they are the only two Bregman loss functions that only depend on the difference (Savage, 1971) or the ratio (Patton, 2011) of the target variable and the forecast.

To illustrate the variety of shapes that Bregman loss functions can take, two parametric families of Bregman loss for variables with support on the real line are presented below. The first was proposed in Gneiting (2011a), and is a family of homogeneous loss functions, where the “shape” parameter determines the degree of homogeneity. We will call this the class of *homogeneous Bregman* loss functions. It is generated by using $\phi(x; k) = |x|^k$ for $k > 1$:

$$L(y, \hat{y}; k) = |y|^k - |\hat{y}|^k - k \operatorname{sgn}(\hat{y}) |\hat{y}|^{k-1} (y - \hat{y}), \quad k > 1 \quad (6)$$

This family nests the squared-error loss function at $k = 2$. (The non-differentiability of ϕ can be ignored if Y_t is continuously distributed, and the absolute value components can be dropped altogether if the target variable is strictly positive, see Patton, 2011).

A second, non-homogeneous, family of Bregman loss can be obtained using $\phi(x; a) = 2a^{-2} \exp\{ax\}$ for $a \neq 0$:

$$L(y, \hat{y}; a) = \frac{2}{a^2} (\exp\{ay\} - \exp\{a\hat{y}\}) - \frac{2}{a} \exp\{a\hat{y}\} (y - \hat{y}), \quad a \neq 0 \quad (7)$$

We will call this the class of *exponential Bregman* loss functions. This family nests the squared-error loss function as $a \rightarrow 0$, and is convenient for obtaining closed-form results when the target variable is Normally distributed, which we exploit below. This loss function has some similarities to the “Linex” loss function, see Varian (1974) and Zellner (1986), in that it involves both linear and exponential terms, however a key difference is that the above family implies that the optimal forecast is the conditional mean, and does not involve higher-order moments.

Figure 1 illustrates the variety of shapes that Bregman loss functions can take and reveals that although all of these loss functions yield the mean as the optimum forecast, their shapes can vary widely: these loss functions can be asymmetric, with either under-predictions or over-predictions being more heavily penalized, and they can be strictly convex or have concave segments. Thus restricting attention to loss functions that generate the mean as the optimum forecast does *not* require imposing symmetry or other assumptions on the loss function. Similarly, in the literature

on economic forecasting under asymmetric loss (see Granger, 1969, Christoffersen and Diebold, 1997, and Patton and Timmermann, 2007, for example), it is generally thought that asymmetric loss functions necessarily lead to optimal forecasts that differ from the conditional mean (they contain an “optimal bias” term). Figure 1 reveals that asymmetric loss functions can indeed still imply the conditional mean as the optimal forecast. (In fact, Savage (1971) shows that of the infinite number of Bregman loss functions, only one is symmetric: the quadratic loss function.)

[INSERT FIGURE 1 ABOUT HERE]

2.2 Forecast comparison in ideal and less-than-ideal forecasting environments

As usual in the forecast comparison literature, I will consider ranking forecasts by their unconditional average loss, a quantity that is estimable, under standard regularity conditions, given a sample of data. (Forecasts themselves, on the other hand, are of course generally based on conditioning information.) For notational simplicity, I assume strict stationarity of the data, but certain forms of heterogeneity can be accommodated by using results for heterogeneous processes, see White (2001) for example. I use t to denote an observation, for example a time period, however the results in this paper are applicable wherever one has repeated observations, for example election forecasting across states, sales forecasting across individual stores, etc.

Firstly, consider a case where forecasters A and B are ranked by mean squared error (MSE)

$$MSE_i \equiv \mathbb{E} \left[\left(Y_t - \hat{Y}_t^i \right)^2 \right], \quad i \in \{A, B\} \quad (8)$$

and we then seek to determine whether

$$MSE_A \lesseqgtr MSE_B \Rightarrow \mathbb{E} \left[L \left(Y_t, \hat{Y}_t^A \right) \right] \lesseqgtr \mathbb{E} \left[L \left(Y_t, \hat{Y}_t^B \right) \right] \quad \forall L \in \mathcal{L}_{Bregman} \quad (9)$$

subject to these expectations existing. The following proposition provides conditions under which the above implication holds. Denote the information sets of forecasters A and B as \mathcal{F}_t^A and \mathcal{F}_t^B .

Assumption 1 *The information sets of the forecasters are nested, so $\mathcal{F}_t^B \subseteq \mathcal{F}_t^A \forall t$ or $\mathcal{F}_t^A \subseteq \mathcal{F}_t^B \forall t$, and do not lead to optimal forecasts that are identical for all t .*

Assumption 2 *If the forecasts are based on models, then the models are free from estimation error.*

Assumption 3 *If the forecasts are based on models, then the models are correctly specified for the statistical functional of interest.*

The above assumptions are presented somewhat generally, as we will refer to them not only in this section on mean forecasting, but also for the analyses of quantile and distribution forecasting below. The second part of Assumption 1 rules out the uninteresting case where two information sets lead to identical forecasts, e.g., they are identical information sets, or where one information set is the union of the other and an information set generated by some random variable that does not lead to a change in the optimal forecast (such as some completely independent random variable). Assumption 3 implies, in this section, that:

$$\exists \theta_{0,i} \in \Theta \text{ s.t. } \mathbb{E} [Y_t | \mathcal{F}_t^i] = m_i (Z_t^i; \theta_{0,i}) \text{ a.s. for some } Z_t^i \in \mathcal{F}_t^i, \text{ for } i \in \{A, B\} \quad (10)$$

where m_i is forecaster i 's prediction model, which has a finite-dimensional parameter θ . The “true” parameter $\theta_{0,i}$ is allowed to vary across i as the conditional mean of Y_t given \mathcal{F}_t^i will generally vary with the information set, \mathcal{F}_t^i . Related, Assumption 2 implies in this section that

$$\hat{Y}_t^i = m_i (Z_t^i; \theta_i^*) \text{ a.s. for all } t = 1, 2, \dots \quad (11)$$

where θ_i^* is some fixed parameter. Part (a) of the proposition below presents a strong, positive, result that holds in the “ideal” forecasting environment. Part (b) shows that a violation of any one of the assumptions in part (a) is sufficient for the positive result to fail to hold.

Proposition 1 *(a) Under Assumptions 1, 2 and 3, the ranking of two forecasts by MSE is sufficient for their ranking by any Bregman loss function.*

(b) If any of Assumptions 1, 2, or 3 fail to hold, then the ranking of two forecasts may be sensitive to the choice of Bregman loss function.

The proof of part (a) is given in the appendix. Of primary interest in this paper is part (b), and we provide analytical examples for this part below.

Under the strong assumptions of comparing only forecasters with nested information sets, and who use only correctly specified models with no estimation error, part (a) shows that the ranking obtained by MSE is sufficient for the ranking by *any* Bregman loss function. This implies that ranking forecasts by a variety of different Bregman loss functions adds no information beyond the MSE ranking. Related to this result, Holzmann and Eulert (2014) show in a general framework that forecasts based on larger information sets lead generally to lower expected loss.

All of the ranking results considered in this paper are in population; in finite evaluation samples rankings of forecasts can switch simply due to sampling variation. If we denote the number of observations available for model estimation and forecast comparison as R and P respectively, then the results here apply for $P \rightarrow \infty$, and when discussing the presence of parameter estimation error (as a violation of Assumption 2) we assume that either R is finite or $R/P \rightarrow 0$ as $R, P \rightarrow \infty$. If instead we consider the case that $P/R \rightarrow 0$ as $R, P \rightarrow \infty$, then we would be in the environment described by Comment 1 to West’s (1996) Theorem 4.1, where parameter estimation error is present but asymptotically negligible. This environment is a generalization of Assumption 2, and all results obtained under Assumption 2 should apply in such an environment.

To verify part (b) of the above proposition, we consider deviations from the three “ideal environment” assumptions used in part (a). Consider the following example: assume that the target variable follows a persistent, but strictly stationary AR(5) process:

$$Y_t = \phi_0 + \phi_1 Y_{t-1} + \dots + \phi_5 Y_{t-5} + \varepsilon_t, \quad \varepsilon_t \sim iid N(0, 1) \tag{12}$$

where $\phi_0 = 1$ and $[\phi_1, \dots, \phi_5] = [0.8, 0.3, -0.5, 0.2, 0.1]$. These parameter values are stylized, but are broadly compatible with estimates for standard macroeconomic time series like US interest rates, see Faust and Wright (2013). We then consider a set of forecasting models. The first three contain no estimation error, and have parameters that are correct given their information sets:

$$\text{AR(1)} \quad \hat{Y}_t = \beta_0 + \beta_1 Y_{t-1} \tag{13}$$

$$\text{AR(2)} \quad \hat{Y}_t = \delta_0 + \delta_1 Y_{t-1} + \delta_2 Y_{t-2} \tag{14}$$

$$\text{AR(5)} \quad \hat{Y}_t = \phi_0 + \phi_1 Y_{t-1} + \dots + \phi_5 Y_{t-5} \tag{15}$$

The first two models use too few lags, while the third model nests the data generating process

and will produce the optimal forecast. The parameters of the first two models are obtained by minimizing the (population) expectation of any Bregman loss function. As each of these models are correctly specified given their (limited) information sets, Proposition 3(a), presented in Section 2.3 below, implies that the optimal parameters are not affected by the specific Bregman loss function used in estimation; I use MSE (making these linear projection coefficients) and present the specific values of the optimal parameters in Appendix SA.1, along with details on the derivation of these parameter values.

In the upper-left panel of Figure 2 I plot the ratio of the expected loss for a given forecast to that for the optimal forecast, as a function of the parameter, a , of the exponential Bregman loss function. (Due to the exponential function, this loss function can lead to large numerical values, which can lead to computational issues in standard software. These can be overcome by simply scaling by some strictly positive value, e.g., the expected loss for the optimal forecast, if available, or some other value.) We see in that panel that the rankings are as expected: the AR(1) model has higher average loss than then AR(2), which in turn has higher average loss than the AR(5). These rankings hold for all values of a , consistent with part (a) of Proposition 1. More generally, the ranking method of Ehm *et al.* (2016) could be applied, and would show that these rankings hold for any Bregman loss function, not only those in the exponential Bregman family.

Now consider comparing two misspecified models. The first is a simple random walk forecast, and the second is a “two-period average” forecast:

$$\text{Random Walk } \hat{Y}_t = Y_{t-1} \tag{16}$$

$$\text{Two-period Average } \hat{Y}_t = \frac{1}{2} (Y_{t-1} + Y_{t-2}) \tag{17}$$

Neither of these forecasts has any estimation error and their information sets are nested, but both models are misspecified. The lower-left panel of Figure 2 compares the average losses for these two forecasts, and we observe that the Random Walk provides the better approximation when the exponential Bregman loss function parameter is near zero, but the Two-period Average forecast is preferred when the parameter is further from zero.

Now we consider the impact of parameter estimation error. Consider the feasible versions of

the AR(2) and AR(5) forecasts, with the parameters are estimated by OLS using a rolling window of 36 observations, corresponding to three years of monthly data:

$$\widehat{\text{AR}}(2) \quad \hat{Y}_t = \hat{\delta}_{0,t} + \hat{\delta}_{1,t}Y_{t-1} + \hat{\delta}_{2,t}Y_{t-2} \quad (18)$$

$$\widehat{\text{AR}}(5) \quad \hat{Y}_t = \hat{\phi}_{0,t} + \hat{\phi}_{1,t}Y_{t-1} + \dots + \hat{\phi}_{5,t}Y_{t-5} \quad (19)$$

We compare $\widehat{\text{AR}}(2)$ and $\widehat{\text{AR}}(5)$ to see whether any trade-off exists between goodness of fit and estimation error: $\widehat{\text{AR}}(5)$ is correctly specified, but requires the estimation of three more parameters; $\widehat{\text{AR}}(2)$ excludes three useful lags, but is less affected by estimation error. Analytical results for the finite-sample estimation error in misspecified AR(p) models are not available, and so we use 10,000 simulated values to obtain the average losses for these two models. The results are presented in the upper-right panel of Figure 2. We see that the expected loss of $\widehat{\text{AR}}(5)$ is below that of $\widehat{\text{AR}}(2)$ for values of the exponential Bregman parameter near zero, while the ranking reverses when the parameter is greater than approximately 0.4 in absolute value. Thus, there is indeed a trade-off between goodness-of-fit and estimation error, and the ranking switches as the loss function parameter changes. This reversal of ranking is not possible in the “ideal environment” case.

Finally, we seek to show that relaxing only Assumption 1 (nested information sets) can lead to a sensitivity in the ranking of two forecasts. For reasons explained below, consider a different data generating process, where the target variable is affected by two independent Bernoulli shocks, X_t and W_t :

$$Y_t = X_t\mu_L + (1 - X_t)\mu_H + W_t\mu_C + (1 - W_t)\mu_M + Z_t \quad (20)$$

where $X_t \sim iid \text{Bernoulli}(p)$, $W_t \sim iid \text{Bernoulli}(q)$, $Z_t \sim iid N(0, 1)$

Forecaster X has access to a “local variation” signal X_t that is regular ($p = 0.5$) but not very strong ($\mu_L = -1$, $\mu_H = 1$), while Forecaster W has access to a “crisis” signal W_t that is irregular ($q = 0.05$) but large when it arrives ($\mu_C = -5$, $\mu_M = 0$). If both forecasters optimally use their (non-overlapping) information sets, then their forecasts are:

$$\hat{Y}_t^X = q\mu_C + (1 - q)\mu_M + \mu_H + (\mu_L - \mu_H)X_t \quad (21)$$

$$\hat{Y}_t^W = p\mu_L + (1 - p)\mu_H + \mu_M + (\mu_C - \mu_M)W_t$$

The lower-right panel of Figure 2 shows that the “crisis” forecaster is preferred for exponential Bregman parameter values less than zero, while the “local variation” forecaster is preferred for larger parameter values.

We have thus demonstrated that relaxing *any one* of the three “ideal environment” assumptions in part (a) of Proposition 1 can lead to sensitivity of forecast rankings to the choice of Bregman loss function. Thus, rather than merely specifying the target functional to be the mean, which narrows the set of relevant loss functions only to the class of Bregman loss functions, forecast consumers or survey designers should specify the *specific* Bregman loss function that will be used to evaluate forecasts. In the next section we consider how this information may be used by forecast producers to better estimate the parameters of their forecasting models.

It should be noted that it may be possible to partially relax Assumptions 1–3 in Proposition 1, or to place other restrictions on the problem, and retain (some, possibly partial) robustness of the ranking of forecasts to the choice of Bregman loss function. One example is when the competing forecasts are correct given their (possibly limited) information sets, free from estimation error, and the target variable and the forecasts are Normally distributed. In this case the following proposition shows we can omit the assumption of nested information sets and retain robustness of rankings for any exponential Bregman loss function. (This explains the need for an alternative DGP in demonstrating sensitivity to non-nested information sets.)

Proposition 2 *If (i) $Y_t \sim N(\mu, \sigma^2)$, (ii) $\hat{Y}_t^i \sim N(\mu, \omega_i^2)$ for $i \in \{A, B\}$, and (iii) $\mathbb{E}[Y_t | \hat{Y}_t^i] = \hat{Y}_t^i$ for $i \in \{A, B\}$, then*

$$MSE_A \lesseqgtr MSE_B \Rightarrow \mathbb{E} \left[L \left(Y_t, \hat{Y}_t^A \right) \right] \lesseqgtr \mathbb{E} \left[L \left(Y_t, \hat{Y}_t^B \right) \right] \quad \forall L \in \mathcal{L}_{Exp-Bregman}$$

Other special cases of robustness may be arise if, for example, the form of the model misspecification was known, or if the target variable has a particularly simple structure (e.g., a binary random variable, see Elliott *et al.* (2016) for example). I do not pursue further special cases here.

2.3 Optimal approximations from a possibly misspecified model

In this section we consider the implications of model misspecification for the *producers* of forecasts. Consider the problem of calibrating a parametric forecasting model to generate the best prediction. If the model is correctly specified, then part (a) of Proposition 3 below shows that minimizing the expected loss under any Bregman loss function will yield a consistent estimator of the model's parameters. We contrast this robust outcome with the sensitivity to the choice of loss function that arises under model misspecification in part (b). Elliott *et al.* (2016) provide several useful related results on this problem when both the target variable and the forecast are binary. They show that even in their relatively tractable case, the presence of model misspecification generally leads to sensitivity of estimated parameters to the choice of (consistent) loss function.

Proposition 3 Denote the model for $\mathbb{E}[Y_t|\mathcal{F}_t]$ as $m(Z_t; \theta)$ where $Z_t \in \mathcal{F}_t$ and $\theta \in \Theta \subseteq \mathbb{R}^p$, $p < \infty$.

Define

$$\theta_\phi^* \equiv \arg \min_{\theta \in \Theta} \mathbb{E}[L(Y_t, m(Z_t; \theta); \phi)] \quad (22)$$

where L is a Bregman loss function characterized by the convex function ϕ . Assume (i) $\partial m(Z_t; \theta) / \partial \theta \neq 0$ a.s. $\forall \theta \in \Theta$ for both (a) and (b) below.

(a) Assume (ii) $\exists! \theta_0 \in \Theta$ s.t. $\mathbb{E}[Y_t|\mathcal{F}_t] = m(Z_t; \theta_0)$ a.s., then $\theta_\phi^* = \theta_0 \forall \phi$.

(b) Assume (ii') $\nexists \theta_0 \in \Theta$ s.t. $\mathbb{E}[Y_t|\mathcal{F}_t] = m(Z_t; \theta_0)$ a.s., then θ_ϕ^* may vary with ϕ .

Assumption (i) in the above proposition is required for identification, imposing that the model is sensitive to changes in the parameter θ . Assumption (ii) is a standard definition of a correctly specified parametric model, and ensures global identification of θ_0 , while Assumption (ii') is a standard definition of a misspecified parametric model.

The proof of part (a) is presented in the appendix. This result is related to the theory for quasi maximum likelihood estimation, see Gouriéroux, *et al.* (1984) and White (1994), for example.

To verify part (b) consider the following illustrative example, where the DGP is:

$$\begin{aligned} Y_t &= X_t^2 + \varepsilon_t, \quad \varepsilon_t \perp\!\!\!\perp X_s \quad \forall t, s \\ X_t &\sim iid N(\mu, \sigma^2), \quad \varepsilon_t \sim iid N(0, 1) \end{aligned} \quad (23)$$

but the forecaster mistakenly assumes the predictor variable enters the model linearly:

$$Y_t = \alpha + \beta X_t + e_t \tag{24}$$

To obtain analytical results to illustrate the main ideas, consider a forecaster using the exponential Bregman loss function defined in equation (7), with parameter a . Using results for functions of Normal random variables (see Appendix SA.1 for details) we can analytically derive the optimal linear model parameters $[\alpha, \beta]$ as a function of a , subject to the condition that $a \neq (2\sigma^2)^{-1}$:

$$\hat{\alpha}_a^* = \sigma^2 - \frac{\mu^2}{(1 - 2a\sigma^2)^2} \quad , \quad \hat{\beta}_a^* = \frac{2\mu}{1 - 2a\sigma^2} \tag{25}$$

This simple example reveals three important features of the problem of loss function-based parameter estimation in the presence of model misspecification. Firstly, the loss function shape parameter does not always affect the optimal model parameters. In this example, if $X \sim N(0, \sigma^2)$, then $(\hat{\alpha}_a^*, \hat{\beta}_a^*) = (\sigma^2, 0)$ for all values of the loss function parameter. Second, identification issues can arise even when the model appears to be *prima facie* well identified. In this example, the estimation problem is not identified at $a = (2\sigma^2)^{-1}$. Issues of identification when estimating under the “relevant” loss function have been previously documented, see Weiss (1996) and Skouras (2007).

Finally, when $\mu \neq 0$, the optimal model parameters will vary with the loss function parameter, and thus the loss function used in estimation will affect the optimal approximation. Figure 3 illustrates this point, presenting the optimal linear approximations for three choices of exponential Bregman parameter, when $\mu = \sigma^2 = 1$. The approximation yielded by OLS regression is obtained when $a = 0$. If we consider a loss function that places more (less) weight on errors that occur for low values of the forecast, $a = -0.5$ ($a = 0.25$) the line flattens (steepens), and Figure 3 shows that this yields a better fit for the left (right) side of the distribution of the predictor variable.

[INSERT FIGURE 3 ABOUT HERE]

The above results motivate declaring the specific loss function that will be used to evaluate forecasts, so that survey respondents can optimize their (potentially misspecified) models taking the relevant loss function into account. It is important to note, however, that it is not always

the case that optimizing the model using the relevant loss function is optimal in finite samples: there is a trade-off between bias in the estimated parameters (computed relative to the probability limits of the parameter estimates obtained using the relevant loss function) and variance (parameter estimation error). It is possible that an efficient (low variance) but biased estimation method could out-perform a less efficient but unbiased estimation method in finite samples. This is related to work on estimation under the “relevant cost function,” see Weiss (1996), Christoffersen and Jacobs (2004), Skouras (2007), Hansen and Dumitrescu (2016) and Elliott *et al.* (2016) for example.

2.4 Comparing quantile forecasts

This section presents results for quantile forecasts that correspond to those above for mean forecasts. The corresponding result for the necessity and sufficiency of Bregman loss for mean forecasts is presented in Saerens (2000), see also Komunjer (2005), Gneiting (2011b) and Thomson (1979): the class of loss functions that is necessary and sufficient for quantile forecasts is called the “generalized piecewise linear” (GPL) class, denoted \mathcal{L}_{GPL}^α :

$$L(y, \hat{y}; \alpha) = (\mathbf{1}\{y \leq \hat{y}\} - \alpha)(g(\hat{y}) - g(y)) \quad (26)$$

where g is a strictly increasing function, and $\alpha \in (0, 1)$ indicates the quantile of interest. A prominent example of a GPL loss function is the “Lin-Lin” (or “tick”) loss function, which is obtained when g is the identity function:

$$L(y, \hat{y}; \alpha) = (\mathbf{1}\{y \leq \hat{y}\} - \alpha)(\hat{y} - y) \quad (27)$$

and which nests absolute error (up to scale) when $\alpha = 1/2$. However, there are clearly an infinite number of loss functions that are consistent for the α quantile. The following is a homogeneous parametric GPL family of loss functions (for variables with support on the real line) related to that proposed by Gneiting (2011b):

$$L(y, \hat{y}; \alpha, b) = (\mathbf{1}\{y \leq \hat{y}\} - \alpha) \left(\text{sgn}(\hat{y}) |\hat{y}|^b - \text{sgn}(y) |y|^b \right) / b, \quad b > 0 \quad (28)$$

Plotting some elements of the homogeneous GPL loss function family (i.e., different choices of b) reveals that their shapes can vary substantially.

When the loss function belongs to the GPL family, the optimal forecast satisfies

$$\alpha = \mathbb{E} \left[\mathbf{1} \left\{ Y_t \leq \hat{Y}_t^* \right\} | \mathcal{F}_t \right] \equiv F_t \left(\hat{Y}_t^* \right) \quad (29)$$

where $Y_t | \mathcal{F}_t \sim F_t$, and if the conditional distribution function is strictly increasing, then $\hat{Y}_t^* = F_t^{-1}(\alpha | \mathcal{F}_t)$. Given its prominence in econometric work, we now seek to determine whether the ranking of two forecasts by Lin-Lin loss is sufficient for their ranking by any GPL loss function (with the same α). That is, whether

$$LinLin_A^\alpha \preceq LinLin_B^\alpha \Rightarrow \mathbb{E} \left[L \left(Y_t, \hat{Y}_t^A \right) \right] \preceq \mathbb{E} \left[L \left(Y_t, \hat{Y}_t^B \right) \right] \quad \forall L \in \mathcal{L}_{GPL}^\alpha \quad (30)$$

subject to these expectations existing. Under the analogous conditions to those for the conditional mean, a sufficiency result obtains.

Proposition 4 (a) *Under Assmptions 1, 2 and 3, the ranking of these two forecasts by expected Lin-Lin loss is sufficient for their ranking by any \mathcal{L}_{GPL}^α loss function.*

(b) *If any of Assumptions 1, 2, or 3 fail to hold, then the ranking of these two forecasts may be sensitive to the choice of \mathcal{L}_{GPL}^α loss function.*

As in the conditional mean case, a violation of any of Assumptions 1, 2, or 3 is sufficient to induce sensitivity to the choice of consistent loss function. A proof of part (a) and analytical examples establishing part (b) are presented in the supplemental appendix. An example based on a realistic simulation design is given in Section 3 below.

2.5 Mean forecasts of symmetric random variables

We next consider a case where some additional information about the target variable is assumed to be known. A leading example in economic forecasting is when the target variable is assumed to be symmetrically distributed. In the following proposition we show that when this assumption holds, the class of loss functions that leads to the forecasters revealing their conditional mean is *even larger* than in the general case in Section 2.1: it is the convex combination of the Bregman and the GPL^{1/2} class of loss functions. The second and third parts present results on ranking forecasters

when the “ideal environment” assumptions hold, or fail to hold. These results suggest that it is even more important to declare which specific loss function will be used to rank the forecasts in such applications, as the set of loss functions that might be employed by survey respondents is even larger than in either the mean (Bregman) or median (GPL^{1/2}) forecasting cases.

Proposition 5 *Assume that $Y_t|\mathcal{F}_{t-1} \sim F_t^*$, a symmetric continuous distribution with finite second moments. Then,*

(a) *Any convex combination of a Bregman and a GPL^{1/2} loss function, $\mathcal{L}_{Breg \times GPL} \equiv \lambda \mathcal{L}_{Bregman} + (1 - \lambda) \mathcal{L}_{GPL}^{1/2}$, $\lambda \in [0, 1]$, yields the mean of F_t^* as the optimal forecast.*

(b) *Under Assmptions 1, 2 and 3, the ranking of these forecasts by MSE or MAE is sufficient for their ranking by any $\mathcal{L}_{Breg \times GPL}$ loss function.*

(c) *If any of Assumptions 1, 2, or 3 fail to hold, then the ranking of these two forecasts may be sensitive to the choice of $\mathcal{L}_{Breg \times GPL}$ loss function.*

2.6 Comparing density forecasts

We now consider results corresponding to the mean and quantile cases above for density or distribution forecasts. In this case the central idea is the use of a proper scoring rule. A “scoring rule,” see Gneiting and Ranjan (2011) for example, is a loss function mapping the density or distribution forecast and the realization to a measure of gain/loss. (In density forecasting this is often taken as a gain, but for comparability with the above two sections I will treat it here as a loss, so that lower values are preferred.) A “proper” scoring rule is any scoring rule such that it is minimized in expectation when the distribution forecast is equal to the true distribution. That is, L is proper if

$$\mathbb{E}_F [L(F, Y)] \equiv \int L(F, y) dF(y) \leq \mathbb{E}_F [L(\tilde{F}, Y)] \quad (31)$$

for all distribution functions $F, \tilde{F} \in \mathcal{P}$, where \mathcal{P} is the class of probability measures being considered. (I will use distributions rather than densities for the main results here, so that they are applicable more generally.) Gneiting and Raftery (2007) show that if L is a proper scoring rule then it must be of the form:

$$L(F, y) = \Psi(F) + \Psi^*(F, y) - \int \Psi^*(F, y) dF(y) \quad (32)$$

where Ψ is a strictly convex, real-valued function, and Ψ^* is a subgradient of Ψ at $F \in \mathcal{P}$. I denote the set of proper scoring rules satisfying equation (32) as $\mathcal{L}_{\text{Proper}}$. As an example of a proper scoring rule, consider the “weighted continuous ranked probability score” from Gneiting and Ranjan (2011):

$$wCRPS(F, y; \omega) = \int_{-\infty}^{\infty} \omega(z) (F(z) - \mathbf{1}\{y \leq z\})^2 dz \quad (33)$$

where ω is a strictly positive weight function on \mathbb{R} . (Strict positivity of the weights makes $wCRPS$ strictly proper.) If ω is constant then the above reduces to the (unweighted) CRPS loss function.

Now we seek to determine whether the ranking of two forecasts by two distribution forecasts by any single proper scoring rule is consistent for their ranking by any proper scoring rule.

$$\mathbb{E}[L_i(F_t^A, Y_t)] \lesseqgtr \mathbb{E}[L_i(F_t^B, Y_t)] \Rightarrow \mathbb{E}[L_j(F_t^A, Y_t)] \lesseqgtr \mathbb{E}[L_j(F_t^B, Y_t)] \quad \forall L_j \in \mathcal{L}_{\text{Proper}} \quad (34)$$

Under the analogous conditions to those for the conditional mean and conditional quantile, a sufficiency result obtains.

Proposition 6 (a) *Under Assmptions 1, 2 and 3, the ranking of these two forecasts by any given proper scoring rule is sufficient for their ranking by any other proper scoring rule.*

(b) *If any of Assumptions 1, 2, or 3 fail to hold, then the ranking of these two forecasts may be sensitive to the choice of proper scoring rule.*

As in the conditional mean and quantile cases, a violation of any of Assumptions 1, 2 or 3 is enough to induce sensitivity to the choice of proper scoring rule. A proof of part (a) and analytical examples establishing part (b) are presented in the supplemental appendix. An example based on a realistic simulation design is given in Section 3 below.

3 Simulation-based results for realistic scenarios

Having established the theoretical possibility of ranking sensitivity in Section 2, the objective of this section is to show that such sensitivity is not a knife-edge result or a mathematical curiosity, but rather a problem that may arise in many practical forecasting applications. I consider three realistic forecasting scenarios, all calibrated to standard economic applications, and show that the

presence of model misspecification, estimation error, or nonnested information sets can lead to sensitivity in the ranking of competing forecasts to the choice of consistent or proper loss functions.

For the first example, consider a point forecast based on a Bregman loss function, and so the target functional is the conditional mean. Assume that the data generating process is a stationary AR(5), with a strong degree of persistence, similar to US inflation or long-term bond yields:

$$Y_t = Y_{t-1} - 0.02Y_{t-2} - 0.02Y_{t-3} - 0.01Y_{t-4} - 0.01Y_{t-5} + \varepsilon_t, \quad \varepsilon_t \sim iid N(0, 1) \quad (35)$$

Now consider the comparison of a parsimonious misspecified model with a correctly-specified model that is subject to estimation error. The first forecast is based on a random walk assumption, and the second forecast is based on a correctly-specified AR(5) model with estimated parameters:

$$\hat{Y}_t^A = Y_{t-1} \quad (36)$$

$$\hat{Y}_t^B = \hat{\phi}_{0,t} + \hat{\phi}_{1,t}Y_{t-1} + \hat{\phi}_{2,t}Y_{t-2} + \hat{\phi}_{3,t}Y_{t-3} + \hat{\phi}_{4,t}Y_{t-4} + \hat{\phi}_{5,t}Y_{t-5} \quad (37)$$

where $\hat{\phi}_{j,t}$ is the OLS estimate of ϕ_j based on data from $t-100$ to $t-1$, for $j = 0, 1, \dots, 5$. I simulate this design for 10,000 observations, and report the differences in average losses for a variety of homogeneous and exponential Bregman loss functions in Figure 4. This figure shows that the ranking of these two forecasts is sensitive to the choice of Bregman loss function: Under squared-error loss (corresponding to parameters 2 and 0 respectively for the homogeneous and exponential Bregman loss functions) the average loss difference is negative, indicating that the AR(5) model has larger average loss than the random walk model, and thus the use of a parsimonious misspecified model is preferred to the use of a correctly specified model that is subject to estimation error. The ranking is reversed for homogeneous Bregman loss functions with parameter above about 3.5, and for exponential Bregman loss functions with parameter greater than about 0.5 in absolute value.

[INSERT FIGURE 4 ABOUT HERE]

Next, consider quantile forecasts for a heteroskedastic time series process, designed to mimic daily stock returns. Such data often have some weak first-order autocorrelation, and time-varying volatility that is well-modeled using a GARCH (Bollerslev, 1986) process:

$$Y_t = \mu_t + \sigma_t \varepsilon_t, \quad \varepsilon_t \sim iid N(0, 1)$$

$$\text{where } \mu_t = 0.03 + 0.05Y_{t-1} \tag{38}$$

$$\sigma_t^2 = 0.05 + 0.9\sigma_{t-1}^2 + 0.05\sigma_{t-1}^2\varepsilon_{t-1}^2$$

I compare two forecasts based on non-nested information sets. The first forecast exploits knowledge of the conditional mean, but assumes a constant conditional variance, while the second is the reverse:

$$\hat{Y}_t^A = \mu_t + \bar{\sigma}\Phi^{-1}(\alpha) \tag{39}$$

$$\hat{Y}_t^B = \bar{\mu} + \sigma_t\Phi^{-1}(\alpha)$$

where $\bar{\mu} = E[Y_t]$, $\bar{\sigma}^2 = V[Y_t]$ and Φ is the standard Normal CDF. I consider these forecasts for two quantiles, a tail quantile ($\alpha = 0.05$) and an intermediate quantile between the tail and the center of the distribution ($\alpha = 0.25$). I compare these forecasts using the family of homogeneous GPL loss functions in equation (28), and report the results based on a simulation of 10,000 observations.

In the right panel of Figure 5, where $\alpha = 0.05$, we see that the forecaster who has access to volatility information (Forecaster B) has lower average loss, across all values of the loss function parameter, than the forecaster who has access only to mean information. This is consistent with previous empirical research on the importance of volatility on estimates of tails. However, when looking at an intermediate quantile, $\alpha = 0.25$, we see that the ranking of these forecasts switches: for loss function parameter values less than about one, the forecaster with access to mean information has lower average loss, while for loss function parameter values above one we see the opposite.

[INSERT FIGURE 5 ABOUT HERE]

As a final example, consider the problem of forecasting the distribution of the target variable. I use a GARCH(1,1) specification (Bollerslev, 1986) for the conditional variance, and a left-skewed t distribution (Hansen, 1994) for the standardized residuals, with parameters broadly designed to match daily US stock returns:

$$Y_t = \sigma_t \varepsilon_t, \quad \varepsilon_t \sim iid Skew t(0, 1, 6, -0.25)$$

$$\sigma_t^2 = 0.05 + 0.9\sigma_{t-1}^2 + 0.05\sigma_{t-1}^2\varepsilon_{t-1}^2 \tag{40}$$

The first distribution forecast is based on the Normal distribution, with mean zero and variance estimated using the past 100 observations. This is a parsimonious specification, but imposes an incorrect model for the predictive distribution. The second forecast is based on the empirical distribution function (EDF) of the data over the past 100 observations, which is clearly more flexible than the first, but will inevitably contain more estimation error.

$$\hat{F}_{A,t}(x) = \Phi\left(\frac{x}{\hat{\sigma}_t}\right), \text{ where } \hat{\sigma}_t^2 = \frac{1}{100} \sum_{j=1}^{100} Y_{t-j}^2 \quad (41)$$

$$\hat{F}_{B,t}(x) = \frac{1}{100} \sum_{j=1}^{100} \mathbf{1}\{Y_{t-j} \leq x\} \quad (42)$$

I consider the weighted CRPS scoring rule (wCRPS) from equation (33) where the weights are based on the standard Normal CDF:

$$\omega(z; \lambda) \equiv \lambda \Phi(z) + (1 - \lambda)(1 - \Phi(z)), \quad \lambda \in [0, 1] \quad (43)$$

When $\lambda = 0$, the scoring rule places more weight on the left tail than the right tail, and the opposite occurs for $\lambda = 1$. When $\lambda = 0.5$ the scoring rule weights both tails equally. Since ω is a convex combination of two weight functions (Φ and $1 - \Phi$), the expected wCRPS is linear in λ .

This design is simulated for 10,000 observations, and the differences in average losses are $[0.51, -0.53, -1.62]$ for $\lambda = [0, 0.5, 1]$. Thus the ranking of these two distribution forecasts is sensitive to the choice of (proper) scoring rule: for weights below about 0.25 (i.e., those with a focus on the left tail), we find the EDF is preferred to the Normal distribution, while for weights above 0.25, including the equal-weighted case at 0.5, the Normal distribution is preferred to the EDF. Thus, the additional estimation error in the EDF generally leads to it being beaten by the parsimonious, misspecified, Normal distribution, *unless* the scoring rule places high weight on the left tail, which is long given the left-skew in the true distribution.

4 Empirical illustration: Evaluating forecasts of US inflation

In this section I illustrate the above ideas using survey forecasts of U.S. inflation. Inflation forecasts are central to many important economic decisions, perhaps most notably those of the Federal Open

Markets Committee in their setting of the Federal Funds rate, but also pension funds, insurance companies, and asset markets more broadly. Inflation is also notoriously hard to predict, with many methods failing to beat a simple random walk model, see Faust and Wright (2013) for a recent comprehensive survey.

Firstly, I consider a comparison of the consensus forecast (defined as the cross-respondent median) of CPI inflation from the Survey of Professional Forecasters (available from tinyurl.com/yckzneb9) and the Thomson Reuters/University of Michigan Survey of Consumers (available from tinyurl.com/y8ef5htj), as well as the Federal Reserve staff “Greenbook” forecasts (available at tinyurl.com/y6vquzq2). For this illustration I examine one-year horizon forecasts, which are directly available for the Michigan and Greenbook forecasts, and can be computed using the one-quarter SPF forecasts for horizons 1 to 4. The sample period is 1982Q3 to 2016Q2, a total of 136 observations, except for the Greenbook forecasts which are only available until 2013Q4 (these forecasts are only available to the public with a five-year lag.) As the “actual” series I use the 2016Q4 vintage of CPI data (available at tinyurl.com/y84skovo). A plot of the forecasts and realized inflation series is presented in Figure 6, and summary statistics are presented in Table 1.

[INSERT FIGURE 6 AND TABLE 1 ABOUT HERE]

I also consider a comparison of individual respondents to the Survey of Professional Forecasters. These respondents are identified in the database only by a numerical identifier, and I select Forecasters 20, 506 and 510, as they all have relatively long histories of responses. (I compare individual forecasters for all periods in which both forecasters are present in the database.) Like the consensus forecasts, I also consider the one-year forecasts from the individual respondents.

Given the difficulty in capturing the dynamics of inflation, it is likely that all forecasters are subject to model misspecification. Further, only relatively few observations are available for forecasters to estimate their model, making estimation error a relevant feature of the problem. Moreover, these forecasts are quite possibly based on nonnested information sets, particularly in the comparison of professional forecasters with the Michigan survey of consumers and the Federal Reserve forecasts. Thus the practical issues highlighted in Section 2 are all potentially relevant here.

Figure 7 presents the results of comparisons of these forecasts, for a range of Bregman loss functions. In the left panels I consider homogeneous Bregman loss functions (equation 6) with parameter ranging from 1.1 to 4 (nesting squared-error loss at 2) and in the right panels I consider exponential Bregman loss functions (equation 7) with parameter ranging from -1 to 1 (nesting squared-error loss at 0). In the top panel we see that the sign of the difference in average losses for SPF and Michigan varies with the parameter of the loss function: the SPF forecast has lower average loss for values of the Bregman parameter less than 3.5 and 0.5 in the homogeneous and exponential cases respectively, while the reverse holds true for parameters above these values. (The difference in average loss is slightly below zero for the squared-error loss case.) The loss difference between the SPF and Greenbook forecasts is positive for values of the Bregman parameter greater than 1.5 and -0.2 in the homogeneous and exponential cases respectively, and (slightly) negative for parameter values less than those thresholds. These results indicate that the ranking of these forecasts of inflation depend on whether over-predictions are more or less costly than under-predictions. (For a given value of the loss function parameter, a Diebold-Mariano (1995) test can be implemented to formally test whether the average loss differences are different from zero. Perhaps unsurprisingly, given the relatively short samples available here, in no case is the null rejected.)

In the middle panel I compare SPF forecaster 20 to forecaster 506, and we again see sensitivity to the choice of loss function: for loss functions that penalize under-prediction more than over-prediction (homogeneous Bregman with parameter less than 2.25, and exponential Bregman with parameter less than zero) forecaster 20 is preferred, while when the loss functions penalize over-prediction more than under-prediction the ranking is reversed. In the lower panel we see an example of a robust ranking: Forecaster 506 has larger average loss than Forecaster 510 for all homogeneous and exponential Bregman loss functions considered; in no case does the ranking reverse.

[INSERT FIGURE 7 ABOUT HERE]

5 Conclusion

Using analytical results, realistic simulation designs, and an application to US inflation forecasting, this paper shows that the ranking of competing forecasts can be sensitive to the choice of consistent loss function or scoring rule. In the absence of model misspecification, parameter estimation error and nonnested forecaster information sets, this sensitivity is shown to vanish, but in almost all practical applications at least one of these complications may be a concern. In the presence of these complications, a conclusion of this paper is that declaring the target functional is not generally sufficient to elicit a forecaster’s best (according to a given, consistent, loss function) forecast; rather best practice for point forecasting is to declare the single, specific loss function that will be used to evaluate forecasts, and to make that loss function consistent for the target functional of interest to the forecast consumer. Reacting to this, forecasters may then wish to estimate their predictive models, if a model is being used, based on the loss function that will evaluate their forecast.

Appendix: Proofs

Proof of Proposition 1(a). We will show that under Assumptions (1)–(3), $MSE_B \geq MSE_A \Rightarrow \mathcal{F}_t^B \subseteq \mathcal{F}_t^A \forall t \Rightarrow \mathbb{E} \left[L \left(Y_t, \hat{Y}_t^B \right) \right] \geq \mathbb{E} \left[L \left(Y_t, \hat{Y}_t^A \right) \right] \forall L \in \mathcal{L}_{Bregman}$.

For the first implication: Assume that $\mathcal{F}_t^A \subseteq \mathcal{F}_t^B \forall t$. This implies $\mathbb{E} \left[\left(Y_t - \hat{Y}_t^A \right)^2 | \mathcal{F}_t^B \right] \geq \mathbb{E} \left[\left(Y_t - \hat{Y}_t^B \right)^2 | \mathcal{F}_t^B \right] a.s. \forall t$ since $\hat{Y}_t^A \in \mathcal{F}_t^A \subseteq \mathcal{F}_t^B$. Then $\mathbb{E} \left[\left(Y_t - \hat{Y}_t^A \right)^2 \right] \geq \mathbb{E} \left[\left(Y_t - \hat{Y}_t^B \right)^2 \right]$ by the law of iterated expectations (LIE). The only way that this can also satisfy the first assumption that $MSE_B \geq MSE_A$ is under equality: $MSE_B = MSE_A$. Since $\mathbb{E} \left[\left(Y_t - \hat{Y}_t^A \right)^2 | \mathcal{F}_t^B \right] \geq \mathbb{E} \left[\left(Y_t - \hat{Y}_t^B \right)^2 | \mathcal{F}_t^B \right] a.s. \forall t$, equality of (unconditional) MSEs can only obtain under equality of conditional MSEs at each point in time, i.e. $\mathbb{E} \left[\left(Y_t - \hat{Y}_t^A \right)^2 | \mathcal{F}_t^B \right] = \mathbb{E} \left[\left(Y_t - \hat{Y}_t^B \right)^2 | \mathcal{F}_t^B \right] a.s. \forall t$, which in turn can only hold if $\hat{Y}_t^A = \hat{Y}_t^B a.s. \forall t$, violating the “not identical” part of Assumption (1). Thus we have a contradiction, and so under Assumptions (1)–(3), $MSE_B \geq MSE_A \Rightarrow \mathcal{F}_t^B \subseteq \mathcal{F}_t^A \forall t$.

Now consider the second implication: Let

$$Y_t = \hat{Y}_t^A + \eta_t = \hat{Y}_t^B + \eta_t + \varepsilon_t \quad (44)$$

Then

$$\begin{aligned}\mathbb{E} \left[L \left(Y_t, \hat{Y}_t^A \right) - L \left(Y_t, \hat{Y}_t^B \right) \right] &= \mathbb{E} \left[-\phi \left(\hat{Y}_t^A \right) - \phi' \left(\hat{Y}_t^A \right) \eta_t + \phi \left(\hat{Y}_t^B \right) + \phi' \left(\hat{Y}_t^B \right) (\eta_t + \varepsilon_t) \right] \\ &= \mathbb{E} \left[\phi \left(\hat{Y}_t^B \right) - \phi \left(\hat{Y}_t^A \right) \right]\end{aligned}\quad (45)$$

since $\mathbb{E} \left[\phi' \left(\hat{Y}_t^A \right) \eta_t \right] = \mathbb{E} \left[\phi' \left(\hat{Y}_t^A \right) \mathbb{E} \left[\eta_t | \mathcal{F}_t^A \right] \right]$ by the LIE and $\mathbb{E} \left[\eta_t | \mathcal{F}_t^A \right] = 0$, by Assumptions (2)-(3). Similarly for $\mathbb{E} \left[\phi' \left(\hat{Y}_t^B \right) (\eta_t + \varepsilon_t) \right]$. Next, consider the second-order mean-value expansion:

$$\phi \left(\hat{Y}_t^A \right) = \phi \left(\hat{Y}_t^B \right) - \phi' \left(\hat{Y}_t^B \right) \varepsilon_t + \phi'' \left(\check{Y}_t^A \right) \varepsilon_t^2 \quad (46)$$

where $\check{Y}_t^A = \lambda_t \hat{Y}_t^A + (1 - \lambda_t) \hat{Y}_t^B$, for $\lambda_t \in [0, 1]$. Thus

$$\mathbb{E} \left[L \left(Y_t, \hat{Y}_t^A \right) - L \left(Y_t, \hat{Y}_t^B \right) \right] = \mathbb{E} \left[\phi' \left(\hat{Y}_t^B \right) \varepsilon_t \right] - \mathbb{E} \left[\phi'' \left(\check{Y}_t^A \right) \varepsilon_t^2 \right] \leq 0 \quad (47)$$

since $\mathbb{E} \left[\phi' \left(\hat{Y}_t^B \right) \varepsilon_t \right] = 0$ and ϕ is convex. And so $\mathcal{F}_t^B \subseteq \mathcal{F}_t^A \forall t \Rightarrow \mathbb{E} \left[L \left(Y_t, \hat{Y}_t^B \right) \right] \geq \mathbb{E} \left[L \left(Y_t, \hat{Y}_t^A \right) \right] \forall L \in \mathcal{L}_{Bregman}$. ■

Proof of Proposition 2. First we note that

$$\begin{aligned}\mathbb{E} \left[L \left(Y_t, \hat{Y}_t^i; a \right) \right] &= \frac{2}{a^2} \left(\mathbb{E} \left[\exp \{ a Y_t \} \right] - \mathbb{E} \left[\exp \{ a \hat{Y}_t^i \} \right] \right) \\ &= \frac{2}{a^2} \left(\exp \left\{ \frac{a}{2} (a\sigma^2 + 2\mu) \right\} - \exp \left\{ \frac{a}{2} (a\omega_i^2 + 2\mu) \right\} \right) \\ &\rightarrow \sigma^2 - \omega_i^2 \text{ as } a \rightarrow 0.\end{aligned}$$

where the first equality holds under mean-unbiasedness (assumption (ii)) and the second follows from normality of the target variable and the forecast (assumption (i)). The last line implies that whichever forecast is based on the richest information set, leading to the greatest (optimal) variability in the forecast (ω_i^2), will have the lowest MSE loss. Then note that for non-MSE exponential Bregman loss (i.e., for $a \neq 0$), that if $\mathbb{E} \left[L \left(Y_t, \hat{Y}_t^A; a \right) \right] \geq \mathbb{E} \left[L \left(Y_t, \hat{Y}_t^B; a \right) \right]$, then $\exp \left\{ \frac{a}{2} (a\omega_A^2 + 2\mu) \right\} \leq \exp \left\{ \frac{a}{2} (a\omega_B^2 + 2\mu) \right\}$ and so $\omega_A^2 \leq \omega_B^2$ and thus $MSE_A \geq MSE_B$. The converse holds using the same derivations, proving the proposition. ■

Proof of Proposition 3(a). The first-order condition for the optimization is:

$$\begin{aligned}
0 &= \left. \frac{\partial}{\partial \theta} \mathbb{E} [L(Y_t, m(X_t; \theta); \phi)] \right|_{\theta = \hat{\theta}_\phi^*} \\
&= \mathbb{E} \left[\phi'' \left(m \left(X_t; \hat{\theta}_\phi^* \right) \right) \left(Y_t - m \left(X_t; \hat{\theta}_\phi^* \right) \right) \frac{\partial m \left(X_t; \hat{\theta}_\phi^* \right)}{\partial \theta} \right] \\
&= \mathbb{E} \left[\phi'' \left(m \left(X_t; \hat{\theta}_\phi^* \right) \right) \left(\mathbb{E} [Y_t | \mathcal{F}_t] - m \left(X_t; \hat{\theta}_\phi^* \right) \right) \frac{\partial m \left(X_t; \hat{\theta}_\phi^* \right)}{\partial \theta} \right]
\end{aligned}$$

where the last equality holds by the LIE. Note that the first-order condition is satisfied when $\hat{\theta}_\phi^* = \theta_0$ by assumption (i), and the solution is unique since ϕ is strictly convex and $\partial m / \partial \theta \neq 0$ a.s. by assumption (ii). ■

The proofs of Propositions 4 and 6 are presented in the supplemental appendix.

Proof of Proposition 5. (a) The first-order condition for an optimal forecast based on a convex combination of Bregman and GPL^{1/2} loss is:

$$0 = -\lambda \phi'' \left(\hat{Y}_t^* \right) \left(\mathbb{E}_{t-1} [Y_t] - \hat{Y}_t^* \right) + (1 - \lambda) \left(\mathbb{E}_{t-1} \left[\mathbf{1} \left\{ Y_t \leq \hat{Y}_t^* \right\} \right] - 1/2 \right) g' \left(\hat{Y}_t^* \right)$$

using the assumption that F_t^* is continuous. Then note that $\mathbb{E}_{t-1} \left[\mathbf{1} \left\{ Y_t \leq \hat{y} \right\} \right] \equiv F_t^* \left(\hat{y} \right)$, and recall that F_t^* is symmetric, which implies that $\mathbb{E}_{t-1} [Y_t] = \text{Median}_{t-1} [Y_t]$ and that $F_t^* \left(\mathbb{E}_{t-1} [Y_t] \right) = 1/2$. Thus $\hat{Y}_t^* = \mathbb{E}_{t-1} [Y_t]$ is a solution to the optimization problem, and this solution is unique as ϕ is strictly convex and g is strictly increasing.

(b) From the proofs of Propositions 1(a) and 4(a), we know that under Assumptions (1)–(3) we have $MSE_B \geq MSE_A \Rightarrow \mathcal{F}_t^B \subseteq \mathcal{F}_t^A \forall t$ and $MAE_B \geq MAE_A \Rightarrow \mathcal{F}_t^B \subseteq \mathcal{F}_t^A \forall t$, and that $\mathcal{F}_t^B \subseteq \mathcal{F}_t^A \forall t \Rightarrow \mathbb{E} \left[L \left(Y_t, \hat{Y}_t^B \right) \right] \geq \mathbb{E} \left[L \left(Y_t, \hat{Y}_t^A \right) \right] \forall L \in \left\{ \mathcal{L}_{Bregman}, \mathcal{L}_{GPL}^{1/2} \right\}$. This immediately yields $\mathcal{F}_t^B \subseteq \mathcal{F}_t^A \forall t \Rightarrow \mathbb{E} \left[L \left(Y_t, \hat{Y}_t^B \right) \right] \geq \mathbb{E} \left[L \left(Y_t, \hat{Y}_t^A \right) \right]$ for any $L \in \mathcal{L}_{Breg \times GPL}$ since $\lambda \in [0, 1]$.

(c) The proof of this negative result requires only an example. This can be constructed using methods similar to those for Propositions 1(b) and 4(b), and is omitted in the interest of brevity.

■

References

- [1] Banerjee, A., X. Guo and H. Wang, 2005, On the Optimality of Conditional Expectation as a Bregman Predictor, *IEEE Transactions on Information Theory*, 51(7), 2664-2669.
- [2] Bauwens, L., C. Hafner and S. Laurent, 2012, *Handbook of Volatility Models and their Applications*, John Wiley & Sons.
- [3] Bollerslev, T., 1986, Generalized Autoregressive Conditional Heteroskedasticity, *Journal of Econometrics*, 31, 307-327.
- [4] Bregman, L.M., 1967, The Relaxation Method of Finding the Common Point of Convex Sets and its Application to the Solution of Problems in Convex Programming, *USSR Computational Mathematics and Mathematical Physics*, 7, 200-217.
- [5] Christoffersen, P. and F.X., Diebold, 1997, Optimal Prediction Under Asymmetric Loss, *Econometric Theory*, 13, 808-817.
- [6] Christoffersen, P. and K. Jacobs, 2004, The Importance of the Loss Function in Option Valuation, *Journal of Financial Economics*, 72, 291-318.
- [7] Diebold, F.X. and R.S. Mariano, 1995, Comparing predictive accuracy, *Journal of Business and Economic Statistics*, 13 (3), 253-263.
- [8] Ehm, W., T. Gneiting, A. Jordan, F. and Krüger, F., 2016, Of quantiles and expectiles: consistent scoring functions, Choquet representations and forecast rankings, *Journal of the Royal Statistical Society, Series B*, 78, 505-562. With discussion and rejoinder.
- [9] Elliott, G., D. Ghanem, and F. Krüger, 2016, Forecasting Conditional Probabilities of Binary Outcomes under Misspecification, *Review of Economics and Statistics*, 98(4), 742-755.
- [10] Elliott, G. and A. Timmermann, 2016, *Economic Forecasting*, Princeton University Press.
- [11] Engelberg, J., C.F. Manski, and J. Williams, 2009, Comparing the Point Predictions and Subjective Probability Distributions of Professional Forecasters, *Journal of Business & Economic Statistics*, 27, 30-41.
- [12] Faust, J. and J.H. Wright, 2013, Forecasting Inflation, in G. Elliott and A. Timmermann (eds.) *Handbook of Economic Forecasting, Volume 2*, Springer Verlag.
- [13] Gourieroux, C., A. Monfort, and A. Trognon, 1984, Pseudo Maximum Likelihood Methods: Theory, *Econometrica*, 52(3), 681-700.
- [14] Gneiting, T., 2011a, Making and Evaluating Point Forecasts, *Journal of the American Statistical Association*, 106(494), 746-762.
- [15] Gneiting, T., 2011b, Quantiles as Optimal Point Forecasts, *International Journal of Forecasting*, 27, 197-207.
- [16] Gneiting, T. and A.E. Raftery, 2007, Strictly Proper Scoring Rules, Prediction and Estimation, *Journal of the American Statistical Association*, 102(477), 358-378.

- [17] Gneiting, T. and R. Ranjan, 2011, Comparing Density Forecasts using Threshold- and Quantile-Weighted Scoring Rules, *Journal of Business & Economic Statistics*, 29(3), 411-422.
- [18] Granger, C.W.J., 1969, Prediction with a Generalized Cost of Error Function, *OR*, 20(2), 199-207.
- [19] Hansen, B.E., 1994, Autoregressive Conditional Density Estimation, *International Economic Review*, 35(3), 705-730.
- [20] Hansen, P.R., and E.-I. Dumitrescu, 2016, Parameter Estimation with Out-of-Sample Objective, Working paper, Department of Economics, UNC-Chapel Hill.
- [21] Holzmann, H. and M. Eulert, 2014, The Role of the Information Set for Forecasting—with Applications to Risk Management, *Annals of Applied Statistics*, 8(1), 595-621.
- [22] Koenker, R., V. Chernozhukov, X. He, and L. Peng, 2017, *Handbook of Quantile Regression*, forthcoming, CRC Press.
- [23] Komunjer, I., 2005, Quasi Maximum-Likelihood Estimation for Conditional Quantiles, *Journal of Econometrics*, 128, 137-164.
- [24] Lieli, R.P. and M.B. Stinchcombe, 2013, On the Recoverability of Forecasters' Preferences, *Econometric Theory*, 29, 517-544.
- [25] Lieli, R.P. and M.B. Stinchcombe, 2017, Unrestricted and Controlled Identification of Loss Functions: Possibility and Impossibility Results, Working paper, Department of Economics, Central European University.
- [26] Merkle, E.C., and M. Steyvers, 2013, Choosing a Strictly Proper Scoring Rule, *Decision Analysis*, 10(4), 292-304.
- [27] Patton, A.J., 2011, Volatility Forecast Comparison using Imperfect Volatility Proxies, *Journal of Econometrics*, 160(1), 246-256.
- [28] Patton, A.J. and A. Timmermann, 2007, Properties of Optimal Forecasts under Asymmetric Loss and Nonlinearity, *Journal of Econometrics*, 140(2), 884-918.
- [29] Saerens, M., 2000, Building Cost Functions Minimizing to Some Summary Statistics, *IEEE Transactions on Neural Networks*, 11, 1263-1271.
- [30] Savage, L.J., 1971, Elicitation of Personal Probabilities and Expectations, *Journal of the American Statistical Association*, 66(336), 783-801.
- [31] Skouras, S., 2007, Decisionmetrics: A Decision-Based Approach to Econometric Modelling, *Journal of Econometrics*, 137, 414-40.
- [32] Thomson, W., 1979, Eliciting Production Possibilities from a Well-Informed Manager, *Journal of Economic Theory*, 20, 360-380.
- [33] Varian, H.R., 1974, A Bayesian Approach to Real Estate Assessment, in S. E. Fienberg and A. Zellner (eds.) *Studies in Bayesian Econometrics and Statistics in Honor of Leonard J. Savage*, North-Holland, Amsterdam, 195-208.

- [34] Weiss, A.A., 1996, Estimating Time Series Models using the Relevant Cost Function, *Journal of Applied Econometrics*, 11, 539-560.
- [35] West, K.D., 1996, Asymptotic inference about predictive ability, *Econometrica*, 64, 1067–1084.
- [36] White, H. 1994, *Estimation, Inference and Specification Analysis*, Econometric Society Monographs No. 22, Cambridge University Press.
- [37] White, H., 2001, *Asymptotic Theory for Econometricians*, 2nd Ed., San Diego, Academic Press.
- [38] Zellner, A., 1986, Bayesian Estimation and Prediction using Asymmetric Loss Functions, *Journal of the American Statistical Association*, 81, 446-451.

Table 1: Summary Statistics

	<i>Mean</i>	<i>Std dev</i>	<i>Min</i>	<i>Max</i>
Actual	2.760	1.373	-1.378	6.255
SPF	3.162	1.279	1.565	8.058
Greenbook	2.917	1.357	0.900	6.400
Michigan	3.163	0.689	1.700	6.900
Forecaster 20	3.432	1.453	0.853	12.142
Forecaster 506	1.640	0.508	-0.047	2.597
Forecaster 510	2.261	0.434	1.256	3.136

Notes: This table presents summary statistics on realized annual US CPI inflation and forecasts of this quantity, all measured in percent, over the period 1982Q3 to 2016Q2. All forecasts are for a one-year horizon. “SPF” is the consensus forecast from the Survey of Professional Forecasters, “Greenbook” is the Federal Reserve staff forecasts (ending in 2013Q4), and “Michigan” is from the Thomson Reuters/University of Michigan Survey of Consumers. The last three rows correspond to forecasts from individual respondents to the Survey of Professional Forecasters, across all observations available for each of the respondents.

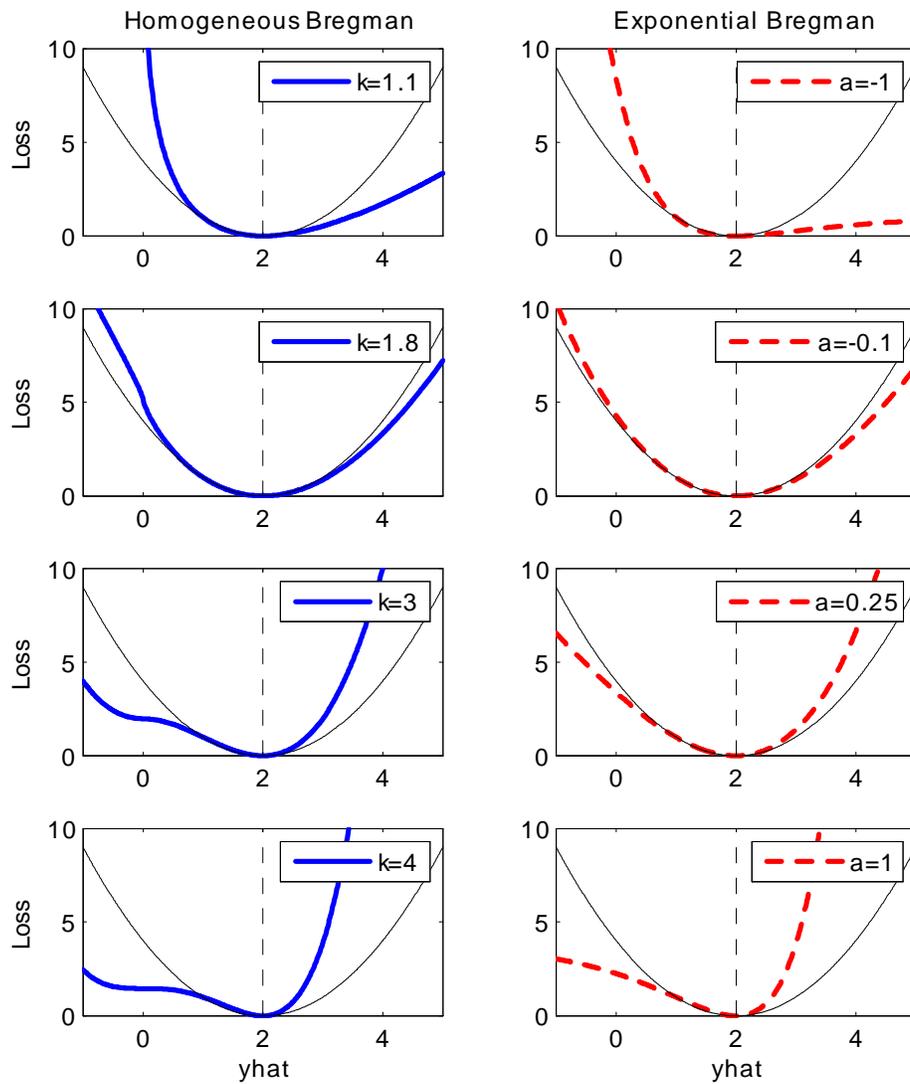


Figure 1: *Various Bregman loss functions. The left column presents four elements of the “homogeneous Bregman” family, and the right column presents four elements of the “exponential Bregman” family. The squared error loss function is also presented in each panel as a thin solid line. In all cases the value for \hat{y} ranges from -1 to 5 , and the value of y is set at 2 .*

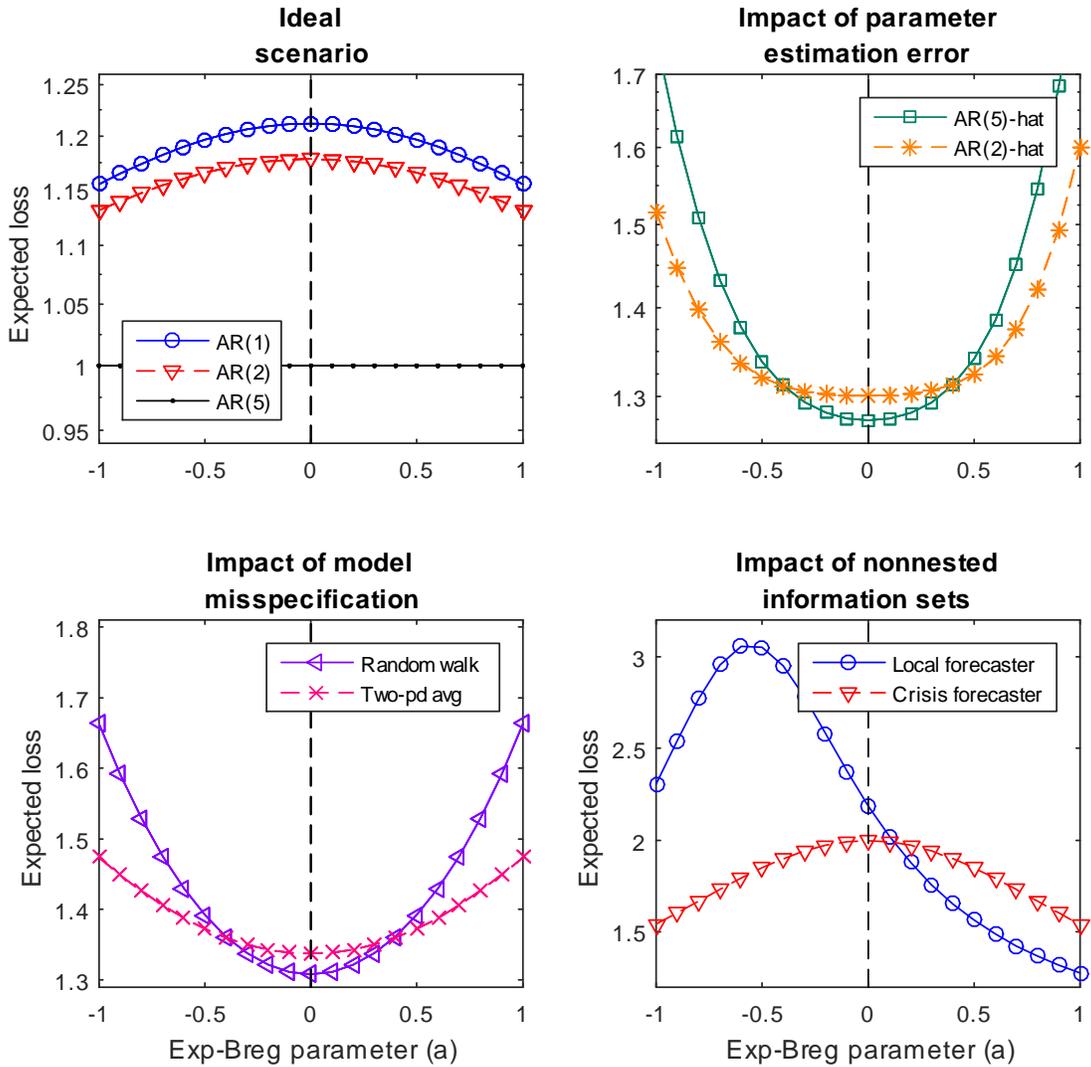


Figure 2: This figure presents the ratio of expected Exponential Bregman loss for a given forecast to that for the optimal forecast, as a function of the Exponential Bregman parameter a .

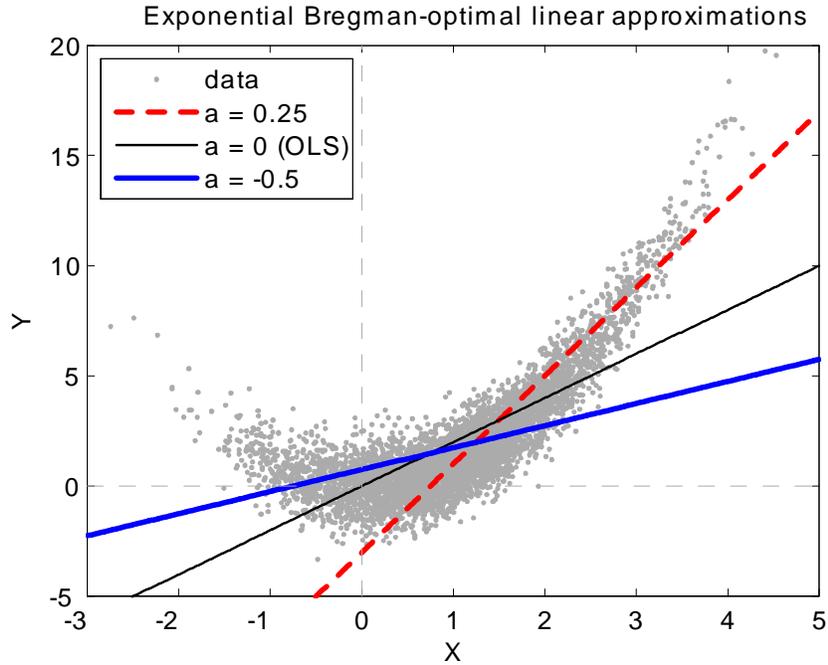


Figure 3: *This figure presents the optimal linear approximations to a nonlinear DGP based on the exponential Bregman loss function for three choices of “shape” parameter; the choice $a=0$ corresponds to quadratic loss, and the fit is the same as that obtained by OLS.*

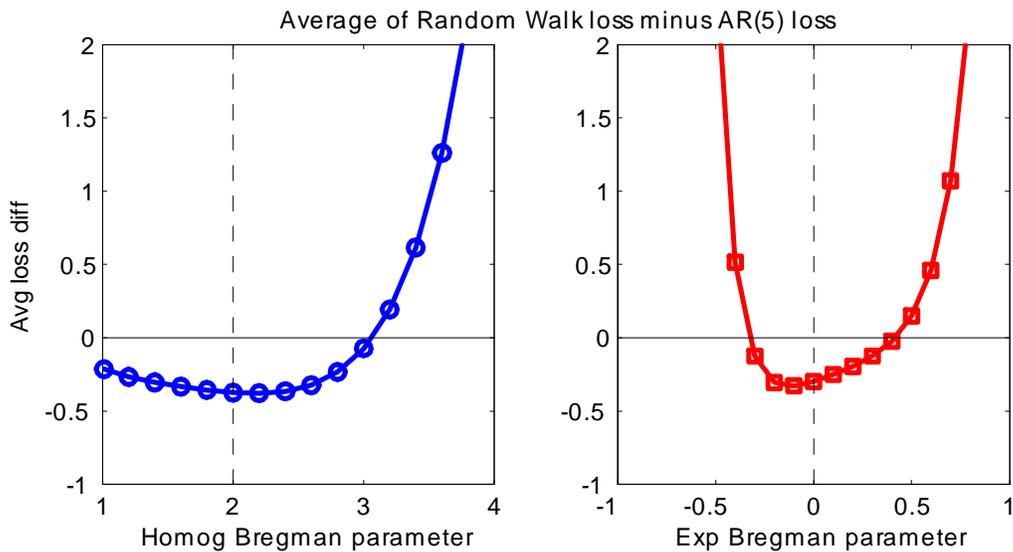


Figure 4: *Average loss from a random walk forecast minus that from an estimated $AR(5)$ forecast, for various homogeneous (left panel) and exponential (right panel) Bregman loss functions.*

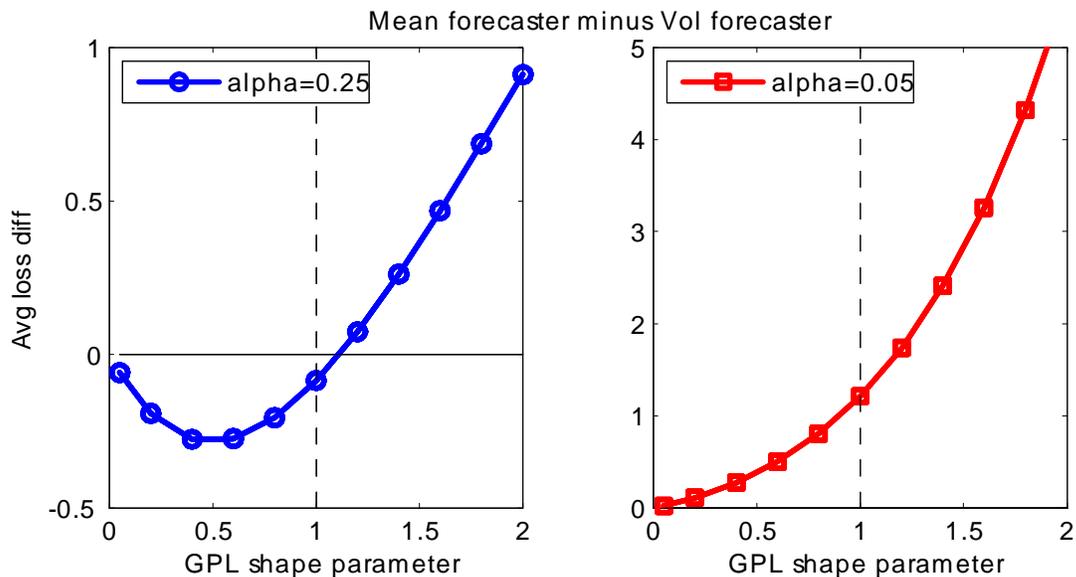


Figure 5: Average loss from a AR-constant volatility forecast minus that from a constant mean-GARCH forecast for various GPL loss functions. (Lin-Lin loss is marked with a vertical line at 1.) The left panel is for the 0.25 quantile, and the right panel is for the 0.05 quantile.

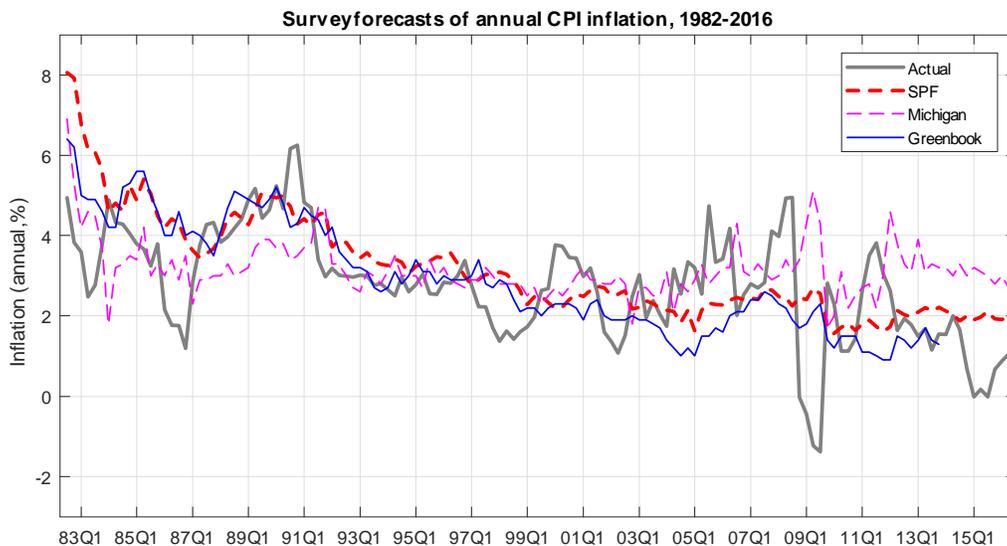


Figure 6: Time series of actual and predicted annual US CPI inflation, updated quarterly, over the period 1982Q3–2016Q2. The inflation forecasts are from the Survey of Professional Forecasters, the Michigan survey of consumers, and the Federal Reserve’s “Greenbook.”

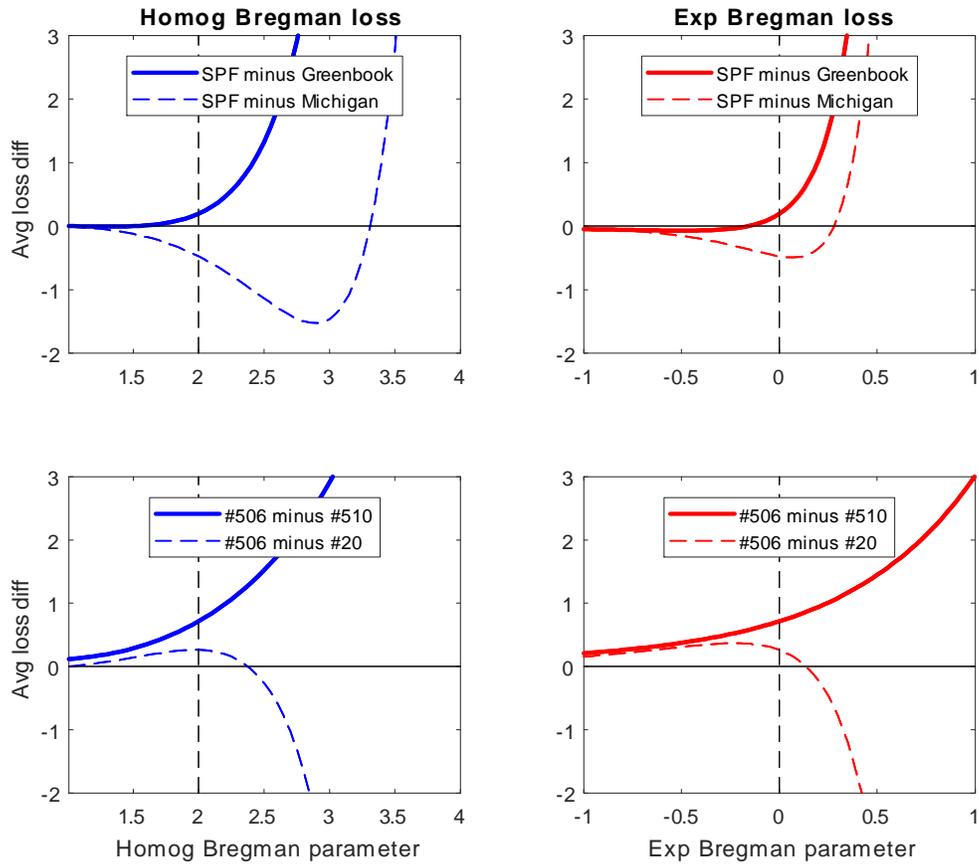


Figure 7: Differences in average losses between two forecasts, for a range of loss function parameters. The “homogeneous Bregman” loss function is in the left column, and the “exponential Bregman” loss function is in the right column. The squared-error loss function is nested at 2 and 0 for these loss functions, and is indicated by a vertical line.