Supplementary Material for "Evaluating infectious disease forecasts with allocation scoring rules"

Aaron Gerding, Nicholas G. Reich, Benjamin Rogers, Evan L. Ray

December 5, 2023

Abstract

We briefly address some technical and methodological points in the main text, referring to the forthcoming ... for more thorough discussion.

| ✓ | From 2.2.1, why are Bayes act scoring rules proper? |
|----------|--|
| | Explain "All proper scoring rules for probabilistic forecasts have an explicit link to a loss function" from discussion. |
| | DGP as optimal for any decision problem, ref Diebold, Gunther, Tay p. 866; and if forecasts are ideal, then forecasts with better information always yield better decisions, ref Holzmann and Eulert, Corr 2. |
| | For 2.2.2, how to get quantile representation of Bayes act using Lagrange multiplier, assuming smooth, never zero densities well behaved at $x = 0$. Work out exponential example. Refer to methods paper for general case. |
| √ | Derivation of quantile scoring rule with quantile as Bayes act for \mathcal{C}/\mathcal{L} problem, assuming never zero densities. |
| | algorithmic details |
| | \square use of distfromq to get from quantiles to distribution functions |
| | \square alloscore |
| | □ implications for propriety. do quantiles elicited by distfromq → alloscore process align with "real" quantiles? the alloscore is proper if distribution functions F are handed to us is it still proper given our algorithm situation? |
| | Descriptions of |
| | \square CRPS as average quantile score across $C \in [0/L]$ decision problems |
| | \square IS as average of two quantile scores with a prob-width penalty |
| | \square WIS as average quantile score across 23 C/L problems. |
| | Sketch of scoring for decision problems involving both cost and constraint. |
| | Derivation of case 2 in formula for Oracle adjustment |

1 Introduction

We briefly address some technical and methodological points in the main text. We begin in section 2 by formalizing the concept of a *shortage* of resources and giving some key results about expected resource shortages under a distribution characterizing uncertainty about (future) levels of resource need. Resource shortages play a central role in the decision-making problems that give rise to the quantile loss and the allocation score, which we discuss in sections 3 and 4 respectively. Section 5 gives details on the numerical methods that we use to calculate allocation scores, including some special

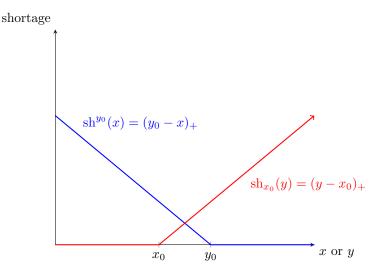


Figure 1: Shortage functions. $\operatorname{sh}^{y_0}(x)$, shown in blue, gives the resource shortage as a function of the level of available resources, x, for a fixed value of resource demand y_0 . $\operatorname{sh}_{x_0}(y)$, shown in red, gives the resource shortage as a function of resource demand, y, for a fixed value of resource supply x_0 .

considerations for settings where forecasts are represented by a finite collection of predictive quantiles, such as the application to forecasts of hospitalizations due to COVID-19 in section 3 of the article.

2 Shortages

We use y to denote the demand or need for resources and x to denote the level of available resources. The resource shortage is the amount by which resource demand exceeds supply. For convenience, we write $u_+ = \max\{0, u\}$, i.e., "the positive part" of u. With this notation, the shortage is written as $(y-x)_+$. To regard shortage as a function of only one variable x or y, with the other being a parameter describing the dependence we can write $(y-x)_+ = \sinh^y(x) = \sinh_x(y)$. Note that $\sinh^y(x)$ and $\sinh_x(y)$ are both convex functions and "mirror" each other:

Let Y be a random variable with distribution F representing the unknown level of resource demand. The random shortage $(Y - x)_+$ can be thought of as either a real-valued random variable $\operatorname{sh}_x(Y)$ for every x, or a function-valued random variable sh^Y whose value for any realization Y = y is a convex function $\operatorname{sh}^y(x)$ of x. In the sections below, we will work with the expected shortage $^1\mathbb{E}_F[(Y - x)_+] = \mathbb{E}_F[\operatorname{sh}^Y](x)$. Assuming that this expected value exists, which is the case as long as the distribution F is well behaved, we can see that $\mathbb{E}_F[\operatorname{sh}^Y](x)$ is also convex (and therefore continuous) in x by integrating the convexity inequality for $\operatorname{sh}^y(x)$ with respect to the probability measure dF(y):

$$\mathbb{E}_{F}[\operatorname{sh}^{Y}](\lambda x_{1} + (1 - \lambda)x_{2}) = \int \operatorname{sh}^{y}(\lambda x_{1} + (1 - \lambda)x_{2})dF(y)$$

$$\leq \int \lambda \operatorname{sh}^{y}(x_{1}) + (1 - \lambda)\operatorname{sh}^{y}(x_{2})dF(y)$$

$$= \lambda \mathbb{E}_{F}[\operatorname{sh}^{Y}](x_{1}) + (1 - \lambda)\mathbb{E}_{F}[\operatorname{sh}^{Y}](x_{2}). \tag{1}$$

¹A more natural sounding term for $(y-x)_+$ might have been *shortfall*. Unfortunately *expected shortfall* has long been used in finance to refer to quantities more closely related to the *conditional* expectation $\mathbb{E}_F[Y-x\mid Y-x\geq 0]=\mathbb{E}_F[(Y-x)_+]/\mathbb{P}_F\{Y\geq x\}.$

Convexity is also shown by directly exhibiting the the left and right derivatives of $\mathbb{E}_F[\operatorname{sh}^Y](x)$:

$$D_{-}\mathbb{E}_{F}[(Y-x)_{+}] = \lim_{h \to 0} \frac{1}{h} \mathbb{E}_{F}[(Y-x)_{+} - (Y-(x-h))_{+}]$$
(2)

$$= \lim_{h \searrow 0} \frac{1}{h} \int_{[x-h,x]} (x-h-y) dF(y) - \lim_{h \searrow 0} \frac{1}{h} \int_{(x,\infty)} h dF(y)$$

$$\tag{3}$$

$$= \lim_{h \searrow 0} \frac{1}{h} \int_{[x-h,x]} -h dF(y) - 1 + F(x) \tag{4}$$

$$= -(F(x) - F(x-)) - 1 + F(x) \quad \left(\text{where } F(x-) := \lim_{t \to x} F(t) \right)$$
 (5)

$$=F(x-)-1\tag{6}$$

$$D_{+} \mathbb{E}_{F}[(Y-x)_{+}] = \lim_{h \searrow 0} \frac{1}{h} \mathbb{E}_{F}[(Y-(x+h))_{+} - (Y-x)_{+}]$$
(7)

$$= \lim_{h \searrow 0} \frac{1}{h} \int_{[x,x+h]} (x-y) dF(y) - \lim_{h \searrow 0} \frac{1}{h} \int_{(x+h,\infty)} h dF(y)$$

$$\tag{8}$$

$$= \lim_{h \searrow 0} \frac{1}{h} \int_{[x,x+h]} 0dF(y) - 1 + F(x) \tag{9}$$

$$= F(x) - 1 \tag{10}$$

where in (4) and (9) we are able to replace the integrands with their values at x because they are bounded over the shrinking regions of integration [x - h, x] and [x, x + h]. Convexity follows since $D_- \mathbb{E}_F[\operatorname{sh}^Y](x) \leq D_+ \mathbb{E}_F[\operatorname{sh}^Y](x)$ by the definition of F(x) and F(x-). This also shows that if F does not have a point mass at x, we have

$$\frac{d}{dx} \mathbb{E}_F[(Y - x)_+] = F(x) - 1, \tag{11}$$

coinciding with the "Leibniz rule" calculation

$$\frac{d}{dx} \mathbb{E}_F[(Y-x)_+] = \frac{d}{dx} \int_x^\infty (y-x) f_Y(y) dy \tag{12}$$

$$= \int_{x}^{\infty} \frac{d}{dx} (y - x) f_{Y}(y) dy - (x - x) f_{Y}(x) = -\int_{x}^{\infty} f_{Y}(y) dy = F(x) - 1.$$
 (13)

which assumes Y has an adequately well-behaved density f_Y .

3 Quantiles and Expected Shortage

We recall how quantiles arise as solutions to a probabilistic decision problem. Let Y be a random variable representing the future level of an undesirable outcome such as severe COVID incidence. Let $x \in \mathbb{R}_+$ be a decision variable representing levels of some costly counter-measure, such as procurement of monoclonal antibody treatments, that can be taken at a cost C > 0 per unit in preparation for Y.² A decision maker must decide on a level x of investment in the counter-measure, and wishes to avoid excesses in either the expediture Cx or the shortage $(y - x)_+$ when Y = y is realized. To formalize the trade-off between these potential excesses we quantify the loss associated with a unit of shortage by a constant L > C (which assumes that the counter-measure has some practical value) and combine the total shortage loss with expenditure into a loss function³

$$l(x,y) = Cx + L(y-x)_{+}.$$

The decision problem is then to select a random future loss l(x, Y) in a way that aligns with the preference that l(x, y) be as low as possible given any realization Y = y.

 $^{^2}$ Quantiles could also be derived for a problem in which x and Y take negative values, corresponding, for instance, to a decision maker that operates in a resource market and a that Y takes negative values when "recoveries" outnumber incidence. But we do not consider such scenarios in this work.

³This does involve a confusing use of the word *loss* to refer to two different quantities, but this seems to be an ingrained and unavoidable habit in the literature.

To give the decision problem more structure we assume the decision maker either knows the distribution F of Y, or wishes to proceed as if a forecast F of Y were true. This gives us what is known in decision theory as a decision problem $under\ risk$ (regarding the future value of Y) as opposed to one $under\ uncertainty$ where both Y as well as F are unknown when the decision is to be made. A principle commonly invoked in this situation⁴ is that the decision maker should or will seek to minimize the expected loss

$$\mathbb{E}_F[l(x,Y)] = Cx + L\mathbb{E}_F[(Y-x)_+]. \tag{14}$$

The expected loss is an affine transformation of the convex expected shortage (c.f. (1)). Therefore $\mathbb{E}_F[l(x,Y)]$ is also is convex and has right and left derivatives $D_{\pm}\mathbb{E}_F[l(x,Y)]$ at every x. Because these derivatives exist everywhere, a necessary condition for x^* to minimize $\mathbb{E}_F[l(x,Y)]$ is that $D_{+}\mathbb{E}_F[l(x^*,Y)] \geq 0$ and $D_{-}\mathbb{E}_F[l(x^*,Y)] \leq 0$, and because of convexity, this condition is also sufficient. From (6) and (10) this means that

$$D_{+}\mathbb{E}_{F}[l(x^{*},Y)] = C + L(F(x^{*}) - 1) \ge 0 \ge D_{-}\mathbb{E}_{F}[l(x^{*},Y)] = C + L(F(x^{*} - 1))$$
(15)

which rearranges with $\alpha = 1 - C/L$ to

$$F(x^*) \ge \alpha \ge F(x^*-). \tag{16}$$

Note that because F(x) and F(x-) are right and left continuous, repectively, and $0 < \alpha < 1$, the set $\{x \mid F(x) \geq \alpha\}$ is closed on the left and the set $\{x \mid \alpha \geq F(x-)\}$ is closed on the right. Therefore, (16) implies that

$$\min\{x \mid F(x) \ge \alpha\} \le x^* \le \max\{x \mid \alpha \ge F(x-)\}. \tag{17}$$

We call $q_{\alpha,F}^- := \min\{x \mid F(x) \ge \alpha\}$ and $q_{\alpha,F}^+ := \max\{x \mid F(x-) \le \alpha\}$ the left and right quantiles of F (for probability level α) and any element $q_{\alpha,F} \in [q_{\alpha,F}^-, q_{\alpha,F}^+]$ a quantile of F. The quantile function for F, which we write as either $Q_F(\alpha)$ or $F^{-1}(\alpha)$, is the set-valued function that maps $\alpha \in (0,1)$ to to the set $[q_{\alpha,F}^-, q_{\alpha,F}^+]$. Thus x^* minimizes the expected loss (14) and gives an optimal solution to the decision problem if and only if $x^* \in Q_F(\alpha)$.

3.1 Quantile functions

For future reference, we record several key properties of quantile functions.

The probability levels $\{\alpha_i\}$ for which $\#Q_F(\alpha_i) > 1$ form a discrete subset of (0,1) and correspond to the non-zero width intervals $[q_{\alpha_i,F}^-, q_{\alpha_i,F}^+]$ where F is constant with values $\{\alpha_i\}$. Conversely, if F is strictly increasing on a (Borel) set $A \subset \mathbb{R}$, then the restriction $Q_F|_{F(A)} = F^{-1}|_{F(A)}$ of Q_F to F(A) is in fact a real-valued function which is left-continuous and provides an inverse to F on F(A). If the the support $\sup(F) = \{x \mid 0 < F(x) < 1\}$ of F is such an F, then extending F^{-1} by $F^{-1}(0) = \inf(\sup(F))$ and $F^{-1}(1) = \sup(\sup(F))$ provides an inverse to F on $F(\mathbb{R}) \subset [0,1]$. If F has a point mass at F is a non-empty set disjoint from F(A), then F takes the constant value F on the closure F is a non-empty set disjoint from F has no discrete component on F, then F is strictly increasing.

Moreover, it can be said that Q_F is increasing as a set-valued function on \mathbb{R} in the generalized sense that $(q_{\alpha,F}-q_{\beta,F})(\alpha-\beta)\geq 0$ whenever $q_{\alpha,F}\in Q_F(\alpha)$ and $q_{\beta,F}\in Q_F(\beta)$, that is, the set graph $(Q_F)=\{(\alpha,q)\mid q\in Q_F(\alpha)\}$ has no downward sloping secants. Conversely, given such an increasing set-valued function Q on (0,1), we can construct a right-continuous function F_Q from \mathbb{R} to [0,1] which will be the cdf of the random variable $Y_Q:=\min(Q(U))$ where $U\sim \mathrm{Unif}[0,1]$.

⁴Note that this priciple might be inappropriate when the decision maker is *risk averse* in some way such as having a preference for random losses with lower variance.

3.2 Opportunity relative to an oracle

Quantiles equivalently arise when the decision problem is defined in terms of the random *opportunity* loss

$$l_o(x,Y) := l(x,Y) - l(Y,Y) = Cx + L(Y-x)_+ - CY$$
(18)

which expresses how much more loss is realized by the decision x than an oracle would have incurred, knowing to invest exactly the future value of Y. The optimal decision for $\mathbb{E}_F[l_o(x,Y)]$ is the same as for $\mathbb{E}_F[l(x,Y)]$ since the term $-C\mathbb{E}_F[Y]$ is constant in x, leading again to the inequalities (15).

Opportunity loss (18) rearranges to

$$l_o(x,Y) = C(x-Y)_+ + (L-C)(Y-X)_+$$
(19)

$$= L(1 - \alpha)(x - Y)_{+} + L\alpha(Y - X)_{+}$$
(20)

$$= L(\alpha - \mathbf{1}\{Y < x\})(Y - x), \tag{21}$$

a form in which it is often called *pinball* loss, despite its graph being an unlikely pinball trajectory for $\alpha \neq 1/2$.

4 Allocation Bayes acts as vectors of marginal quantiles.

Here we show that the Bayes act $x^{F,K}=(x_1^{F,K},\ldots,x_N^{F,K})$ for a forecast F, corresponding to the allocation problem (??) (in Section ?? in the main text) can be represented as a vector of quantiles for the marginal forecast distributions F_i at a single probability level $\tau^{F,K}$, that is, $x_i^{F,K}=q_{F_i,\tau^{F,K}}$. An immediate consequence used in the examples in Section ?? in the main text is that if $F_i=\text{Exp}(1/\sigma_i)$ for all i, then the Bayes act is proportional to $(\sigma_1,\ldots,\sigma_N)$, since $q_{\text{Exp}(1/\sigma),\tau}=-\sigma\log(1-\tau)$.

For an arbitrary allocation vector $x \in \mathbb{R}^N_+$ the expected loss

$$\mathbb{E}_{F}[s_{A}(x,Y)] = \sum_{i=1}^{N} L \cdot \mathbb{E}_{F_{i}}[(Y_{i} - x_{i})_{+}]$$
(22)

is the sum of expected shortages (scaled by L) under the allocations x_i in each location. We therefore have the following necessary condition for $x^* \in \mathbb{R}^N_+$ to be an optimal allocation for $\mathbb{E}_F[s_A(x,Y)]$ under the constraint $\sum_{i=1}^N x_i = K$: if $\delta > 0$ of the x_i^* units of resource allocated to location i are reallocated to location j, expected shortage will increase in location i by at least as much as it decreases in location i. That is,

$$\mathbb{E}_{F_i}[(Y_i - (x_i^* - \delta))_+] - \mathbb{E}_{F_i}[(Y_i - x_i^*)_+] \ge \mathbb{E}_{F_i}[(Y_j - x_j^*)_+] - \mathbb{E}_{F_i}[(Y_j - (x_j^* + \delta))_+]. \tag{23}$$

Since the expected shortages in i and j have right and left derivatives at any x_i and x_j (see Section 2), we can divide (23) by δ and take limits for $\delta \searrow 0$ to get

$$-D_{-}\mathbb{E}_{F}[(Y_{i} - x_{i}^{\star})_{+}] \ge -D_{+}\mathbb{E}_{F}[(Y_{i} - x_{i}^{\star})_{+}]. \tag{24}$$

Note that the minus signs appear because our optimality condition addresses how a *decrease* in resources will *increase* the expected shortage in i and vice versa in j. Scaling by L to match the right and left partial derivatives of $\mathbb{E}_F[s_A(x,Y)]$ and using formulae (6) and (10), (24) becomes

$$L(1 - F_i(x_i^* -)) \ge L(1 - F_i(x_i^*)).$$
 (25)

Inequalities (24) and (25) remain true with i and j reversed. They hold with i = j as well by the definition of $F_i(x_i^*-)$. Therefore, a number λ (a Lagrange multiplier) exists, which is independent of i, such that

$$L(1 - F_i(x_i^* -)) \ge \lambda \ge L(1 - F_i(x_i^*)), \quad \text{for all } i \in 1, \dots, N.$$

That is,

$$F_i(x_i^*) \ge 1 - \lambda/L \ge F_i(x_i^* -),\tag{27}$$

which says (c.f. discussion after (16) and (17)) that x_i^* is a quantile q_{τ,F_i} for $\tau = 1 - \lambda/L$.

The constraint now determines τ (and hence the Bayes act) through $\sum_{i=1}^{N} q_{\tau,F_i} = K$. This equation can be restated as the requirement that

$$K \in TQ_F(\tau) \tag{28}$$

where the set-valued function

$$TQ_F(\tau) := \sum_{i=1}^{N} Q_{F_i}(\tau) = \left[\sum_{i=1}^{N} q_{\tau,F_i}^-, \sum_{i=1}^{N} q_{\tau,F_i}^+ \right]$$
 (29)

is defined using interval addition [a, b] + [c, d] = [a + c, b + d]. (Note that the letter T is being used to connote a "totalling" operation, as it often is in survey sampling literature.)

 TQ_F satisfies the conditions mentioned in section 3.1 for being a quantile function, and so there is a random variable TY_F with cdf $F_T := F_{TQ_F}$. From this perspective, the problem of finding τ becomes the calculation of $F_T(K) = \mathbb{P}(TY_F \leq K)$, making clear the existence of a solution τ^* to (28). This also yields the interesting formal equation for the Bayes act

$$\mathbf{F}(x^{F,K}) = \mathbf{1}_N F_T(K),\tag{30}$$

where $\mathbf{1}_N : \mathbb{R} \to \mathbb{R}^N$ is the linear map that takes a to the N-vector $(a, \dots, a)^T$ and the vector of marginal cdfs $\mathbf{F} := (F_1, \dots, F_N)$ is a map from \mathbb{R}^N to \mathbb{R}^N . With this notation we can write (28) as

$$K \in TQ_F\left(\frac{1}{N}\mathbf{1}^T \mathbf{\tilde{F}}(x^{F,K})\right),$$
 (31)

which leads conceptually to the iterative numerical method of solving for τ and $x^{F,K}$ discussed next in section 5.

Two awkward features of the quantile representation of the Bayes act can arise. First, point masses in the F_i create point masses for F_T which may cause τ^* to not be the unique solution to (28). Secondly, if more than one $Q_{F_i}(\tau^*)$ is a positive-width interval, then the Bayes act will not be unique in these coordinates, and generically not all points in the $Q_{F_i}(\tau^*)$ will be the coordinate of a Bayes act.

It is important to note that λ depends on the forecast F and the constraint level K. Thus while $\lambda = L(1-\tau)$ can be interpreted as a kind of "cost" imposed by the constraint in the allocation problem which is analogous to $C = L(1-\alpha)$ in the the cost-lost problem of section 3, it does not serve to define the allocation loss function in the way that C defines (14). λ is rather a parameter that must be found, given the pair F and K.

5 Numerical computation of allocation scores

Suppose we have established that $\tau^* = F_T(K)$ lies in the interval $I_1 = [\tau_L, \tau_U]$ with $\tau_L < \tau_U$, that is, $K \in [q_{F_T, \tau_L}^-, q_{F_T, \tau_U}^+]$. From section 3.1, we know that the set $TQ_F(\tau_L) \cup TQ_F(\tau_U) \subset [q_{F_T, \tau_L}^-, q_{F_T, \tau_U}^+]$ is arranged in exactly one of the following ways:

- (••) $TQ_F(\tau_L) \cap TQ_F(\tau_U) = \emptyset$ (and $q_{F_T,\tau_L}^+ < q_{F_T,\tau_U}^-$)
- (•) $TQ_F(\tau_L) = \{K\} = TQ_F(\tau_U)$ (a point mass at K)
- (•-) $TQ_F(\tau_L) = \{q_{F_T,\tau_U}^-\} \subsetneq TQ_F(\tau_U)$ (a point mass at q_{F_T,τ_U}^-)
- $(- \bullet) \ TQ_F(\tau_L) \supsetneq \{q^+_{F_T,\tau_L}\} = TQ_F(\tau_U)$ (a point mass at $q^+_{F_T,\tau_U})$

In the case (\bullet), we can immediately take $\tau^* = \tau_U$ as the probability level defining the allocation Bayes act. In the cases (\bullet -), and ($-\bullet$), which imply the presence of a point mass in one or more of the component forecasts adjacent to a region of zero density in all components, we can take $\tau^* = \tau_U$ or τ_L , respectively, as the representing probability level and find Bayes acts within the set $\mathbf{F}^{-1}(\tau^*)$ (by solving the linear program $\sum_{x \in \mathbf{F}^{-1}(\tau^*)} x_i = K$).

Otherwise, by evaluating TQ_F at $\tau_M = \frac{1}{2} (\tau_L + \tau_U)$ we can replace I_1 with the interval

$$I_{2} = \begin{cases} [\tau_{L}, \tau_{M}] & \text{if } K < q_{F_{T}, \tau_{M}}^{-} \\ [\tau_{M}, \tau_{U}] & \text{if } K \ge q_{F_{T}, \tau_{M}}^{-}. \end{cases}$$
(32)

Iterating this process we obtain a sequence $\{I_k\}, k=1,2,\ldots$ of intervals of widths $|I_k|=2^{1-k}|I_1|$ which either terminates at one of the scenarios (\bullet) , $(\bullet-)$, or $(-\bullet)$, or provides infinite sequences $\{\tau_{L,k}\}$ and $\{\tau_{U,k}\}$ converging to τ^* from below and above.

In the generic case of an infinite $\{I_k\}$, we need to define practical stopping criteria for each of three possible behaviours of F_T at K:

- (\mathscr{I}) $\tau^* = F_T(K-)$ and $TQ_F(\tau^*) = \{K\}$ (that is, if F_T is continuous and strictly increasing at K)
- $(_ \bullet -)$ F_T has a point mass at K
- (\bullet) $K \in (q_{F_T,\tau^*}^-, q_{F_T,\tau^*}^+)$ (that is, all component forecasts F_i have zero density at points $x_i \in (q_{F_i,\tau^*}^-, q_{F_i,\tau^*}^+)$ and $\sum x_i = K$).

6 Properties and Properness

For a prediction to be useful, it must **proper**ly describe a **proper**ty.

Expanding on the decision theoretic perspective sketched in Section 3, we can view a loss function as a general tool for formalizing a decision problem that assigns numerical value to the result of taking an action x in preparation for an outcome y. A scoring rule S is a loss function where the action is a probabilistic forecast F of the outcome y (or the statement of F by a forecaster). Just as in Section 3, given an action F, S transforms a random outcome variable Y into a random loss S(F,Y). We refer to the realized loss S(F,y) as the score of F at y, and the process of evaluating $S(F,y_i)$ for a data set $Y = \{y_i\}$ as scoring F against Y. Perhaps the most fundamental example of a scoring rule from a statistical perspective is the logarithmic score $S_{\log}(F,y) = -\log f(y)$ (where f(y) is the density or the mass of F at y if it exists and otherwise the mass), that is, the negative log-likelihood of F interpreted as a parameter for the singleton data set $\{y\}$. This frames a maximum likelihood estimate as the solution to a decision problem following the expected loss minimization principle introduced in Section 3.

Decision theoretically, probabilistic forecasts are a unique kind of action in that they can be used to generate their own(simulated) outcome data, against which they can be scored using S. S therefore commits a probabilistic forecast F to the "self-assessment" $\mathbb{E}[S(F,Y^F)]$, where $Y^F \sim F$ is the random variable defined by sampling from F, as well to an assessment $\mathbb{E}[S(G,Y^F)]$ of any alternative forecast G. For S_{\log} this self-assessment is the Shannon entropy $H(F) = -\int \log(f(x))f(x)dx$ of F, and adding to H(F) the Kullback-Leibler (KL) divergence $D_{KL}(F,G) = -\int \log(g(x)/f(x))f(x)dx$ gives F's assessment $\mathbb{E}[S_{\log}(G,Y^F)]$ of G. (That is, $\mathbb{E}[S_{\log}(G,Y^F)] = H(F) + D_{KL}(F,G)$.) The KL divergence is the degree to which F perceives G as divergent from being able to minimize expected loss in this particular forecaster decision problem (i.e. maximize expected log-likelihood).

A natural consistency criterion for S is that it does not commit F to assessing any other forecast G as being better than F itself, that is, that

$$\mathbb{E}[S(F, Y^F)] \le \mathbb{E}[S(G, Y^F)] \tag{33}$$

for any F, G. Otherwise, the optimal decision for some forecaster would be to state a forecast G other than the forecast F which they believe describes the outcome Y. A scoring rule meeting this criterion

is called *proper*. The inequality can also be written as $\mathbb{E}_F[S(F,Y)] \leq \mathbb{E}_F[S(G,Y)]$ where the subscript specifies the distribution of Y. S is *strictly proper* when this inequality is sharp, in which case the *only* optimal decision for a forecaster is to state the forecast they believe to be true. The logarithmic scoring rule S_{\log} , for example, is strictly proper due to the positivity of KL divergence.

Another definition of S being proper that is often used is that $\mathbb{E}[S(F,Y)]$ is lowest when F is the true distribution of Y. Under a flexible reading, this definition is equivalent to ours, but we find it problematic because it invites the mistaken impression that the properness of a score might depend on the true distribution of what is being forecasted. Whether a score is proper is unrelated to any particular forecast being scored or source of the data being used to score it.

The condition of being proper is quite strong, and naïve means of reducing a forecast distribution to a single number based on an observed outcome y will generally define an improper scoring rule. For example, the probability score $S_{\text{Prob},c}(F,y) := -(F(y+c)-F(y-c))$ and the linear score $S_{\text{Lin}}(F,y) := \lim_{c\to 0} \frac{1}{2c} S_{\text{Prob},c}(F,y) = -f(y)$ (where we assume F has a density f) and are both improper because they commit a general F to assessing an alternative forecast G as being better than F itself whenever G is sufficiently more concentrated than F in the neighborhood $[m_f - c, m_f + c]$ of a mode m_f of f. The same is true when the outcome is discrete (and ordered for S_{Prob}) with index y: for a forecast pmf p(y), moving the mass of p onto a neigborhood of a mode m_p to get a new forecast p_m will improve the expected score $\mathbb{E}[S_{\text{Prob/Lin}}(p_m, Y^p)]$ according to p itself. This creates the classic pathology of probabilistic forecasters having no incentive to express uncertainty. The meteorologist getting paid according to p_{rain} rain $+(1-p_{\text{rain}})(1-\text{rain})$ (the negative of the linear score) will say rain is inevitable or impossible when the chances of rain appear to be 51% or 49%.

Sometimes a scoring rule can be made proper with an added correction that penalizes

But given an auxiliary decision problem $\mathcal{D}_{Aux} = \{\mathcal{X}, \mathcal{Y}, l(x,y)\}$, such as the allocation problem from the main text or the cost-loss problem from Section 3, we can produce a scoring rule S_l that is automatically proper using the Bayes act formalism introduced in Section ??. This works by defining $S_l(F,y)$ as the value $l(x^F,y)$ already given by \mathcal{D}_{Aux} for the Bayes act x^F ; that is, the action which by design will be assessed via the forecast F to have the lowest possible expected loss $\mathbb{E}[l(x^F,Y^F)]$. In statistical decision theory, $\mathbb{E}[l(x^F,Y^F)] = \mathbb{E}_F[l(x^F,Y)]$ is sometimes referred to as the Bayes risk of F relative to \mathcal{D}_{Aux} .

by mapping a forecast F to an the action x^F

The study of elicitability also concerns this construction (see, e.g., Gneiting [2011]), but from a different vantage point than ours. The question there is whether for a given map (called a functional in statistics and a property in computer science) $T(F) = x \in \mathcal{X}$ of distributions into the auxiliary action space, there exists a loss function l such that $T(F) = x^F$ is the Bayes act for l and the associated scoring rule $S_l(F,x)$ is strictly proper. Such an l is said to elicit T. Our focus is rather on how forecast evaluation proceeds via the Bayes act construction for a given loss function of subject matter interest.

References

Morris H DeGroot. Optimal statistical decisions. John Wiley & Sons, 2005.

James O Berger. Statistical decision theory and Bayesian analysis. Springer Science & Business Media, 2013.

Tilmann Gneiting. Making and evaluating point forecasts. *Journal of the American Statistical Association*, 106(494):746–762, 2011.

 $^{^5}$ A clear exposition of basic decision theory using this terminology is Chapter 8 of DeGroot [2005]. Unfortunately other influential sources use "Bayes risk" to refer to a wider variety of quantities. In Berger [2013] for example, the Bayes risk can be both the optimal expected loss for F (see p. 17) or the expected loss for a general x with respect to F (see p. 6).