

PREDICTION IN THE WORST CASE¹

BY DEAN P. FOSTER

University of Chicago

A predictor is a method of estimating the probability of future events over an infinite data sequence. One predictor is as strong as another if for all data sequences the former has at most the mean square error (MSE) of the latter. Given any countable set \mathcal{D} of predictors, we explicitly construct a predictor S that is at least as strong as every element of \mathcal{D} . Finite sample bounds are also given which hold uniformly on the space of all possible data.

1. Introduction and summary. As in Rissanen (1986), the data will be an infinite sequence of binary variables R_i , in other words, $R = \{R_i, \text{ for } i = 1, 2, \dots\}$ is a $\{0, 1\}$ -valued stochastic process. Ω is the set of infinite binary sequences. Let \mathcal{F}_n be the natural σ -field induced by R_1, R_2, \dots, R_n . Call M a predictor if at time $n - 1$ it outputs a number $M(n)$ which estimates the probability $P(R_n = 1 | \mathcal{F}_{n-1})$. Technically, M is a predictor if and only if it is an \mathcal{F}_{n-1} -adapted stochastic process taking values in the interval $[0, 1]$. This model of probability forecasting is discussed in Dawid (1984, 1986).

Consider two predictors: A and B . Which is better? Define

$$\bar{Q}_R(A, n) = \frac{1}{n} \sum_{i=1}^n (A(i) - R_i)^2,$$

the mean square error between prediction and observation—called the Brier score in meteorology [Brier (1950)]. We can compare predictors A and B by computing $d_R(A, B)$,

$$d_R(A, B) = \liminf_{n \rightarrow \infty} (\bar{Q}_R(A, n) - \bar{Q}_R(B, n)).$$

Note that $d_R(\cdot, \cdot)$ is not a metric, nor is it antisymmetric. Since \bar{Q}_R is bounded between zero and one, $d_R(\cdot, \cdot)$ is finite. If $d_R(A, B) > 0$, then in the limit B forecasts better than A on the data sequence R .

DEFINITION. B is as strong as A if for all sequences R in Ω , $d_R(A, B) \geq 0$.

DEFINITION. B is stronger than A if it is as strong as A and there exists a sequence R in Ω such that $d_R(A, B) > 0$.

The results of this paper concern relationships among predictors, such as strong or stronger, which hold for all sample paths. Although not explicitly

Received October 1987; revised April 1990.

¹Research supported by the Office of Naval Research under Contract FAS-01-5-28870 at the University of Maryland.

AMS 1980 subject classifications. Primary 62M20; secondary 62A99.

Key words and phrases. Comparing forecasts, worst-case behavior, mean square error.

mentioned, such properties will hold a fortiori almost surely or in expectation with respect to any measure \mathcal{P} on Ω . For example, if B is as strong as A , then for all measures \mathcal{P} ,

$$\liminf_{n \rightarrow \infty} E(\bar{Q}_R(A, n)) - E(\bar{Q}_R(B, n)) \geq 0.$$

The goal of this paper is to propose and justify a new method for combining predictors. More precisely, consider two predictors A and B , neither of which is as strong as the other. This new method will yield a predictor S which is as strong as both A and B . In other words, S combines the best properties of both A and B ; if A predicts well on some sequence of R_i , then so does S .

Empirical properties of various other methods of combining forecasts have been discussed in Newbold and Granger (1974) and Gupta and Wilton (1987). A concept similar to strength which deals with portfolio selection in finance is due to Cover (1991).

2. Exploring strength.

EXAMPLE. Define two predictors: $ZERO(n) \equiv 0$ and $ONE(n) \equiv 1$. Then, ZERO is not as strong as ONE, nor is ONE as strong as ZERO. This is seen by considering the two data sequences $R_i = 1$ for all i and the sequence $R'_i = 0$ for all i . Then, $d_R(ONE, ZERO) = -1$, $d_R(ZERO, ONE) = -1$.

This situation seems reasonable—each of these two predictors has some data that it does better on. Is stronger a vacuous definition? No; the following lemma gives a general construction.

LEMMA 1. *For any arbitrary predictor M , there exists a predictor S which is stronger than M .*

PROOF. We need to find two things to prove this lemma: a predictor S which is as strong as M and a point ω in Ω for which this predictor does better than M . Find a point ω in Ω with the following property:

$$\omega_i = \begin{cases} 1, & \text{if } M(i)(\omega) \leq 0.5, \\ 0, & \text{if } M(i)(\omega) > 0.5. \end{cases}$$

This point can be recursively constructed because $M(i)$ is \mathcal{F}_{i-1} -measurable. Note that $\bar{Q}_\omega(M, n) \geq 0.25$ for all n : Now let

$$S(i)(R) = \begin{cases} \omega_i, & \text{if } R_j = \omega_j \text{ for all } j < i, \\ M(i), & \text{otherwise.} \end{cases}$$

Thus, for all points $R \neq \omega$ in Ω , $d_R(M, S) = 0$. Also, $d_\omega(M, S) \geq 0.25$. Thus, S is stronger than M . \square

One would expect from the name that strength is transitive. This is so. It follows from superadditivity of limit inferior: $d_R(C, A) \geq d_R(C, B) + d_R(B, A) \geq 0$.

3. Results. Suppose we have a finite set \mathcal{D} of predictors. If a strongest element of \mathcal{D} does not exist, then we cannot decide which predictor from \mathcal{D} to use. Thus, the search is on for finding a predictor $S_{\mathcal{D}}$ which is as strong as every element of \mathcal{D} . The subscript \mathcal{D} in $S_{\mathcal{D}}$ stresses that S depends on which set is being dominated.

Call the elements of \mathcal{D} : $f_1, f_2, f_3, \dots, f_d$. For instance, the prediction made by the second predictor in \mathcal{D} at time n is $f_2(n)$. Define the convex hull of \mathcal{D} to be the set of predictors created by convex weightings of predictors in \mathcal{D} .

DEFINITION. Convex hull of \mathcal{D} is $\{M | \dot{M}(i) = w_1 f_1(i) + w_2 f_2(i) + \dots + w_d f_d(i) \text{ for } w_i \geq 0 \text{ and } w_1 + w_2 + \dots + w_d = 1\}$.

In vector notation, $M = w^T f$, where $w = [w_1, w_2, \dots, w_d]^T$ and $f = [f_1, f_2, \dots, f_d]^T$. At time n , we write the prediction made by M as $M(n) = w^T f(n)$.

DEFINITION. Let B represent the simplex

$$B = \left\{ w \in \mathbb{R}^d \mid \sum_{i=1}^d w_i = 1, w_i \geq 0 \right\}.$$

The convex hull of \mathcal{D} can be written as $\{M | M = w^T f, w \in B\}$.

We now state the main theorem.

THEOREM 1. For \mathcal{D} a collection of d predictors, let

$$S_{\mathcal{D}}(n) = \hat{w}(n-1)^T f(n),$$

where $\hat{w}(n-1)$ minimizes $(n-1)\bar{Q}(w^T f, n-1) + w^T w$ over w in B . Then $S_{\mathcal{D}}$ is as strong as every element of the convex hull of \mathcal{D} . In particular: $\bar{Q}_R(S_{\mathcal{D}}, n) - \bar{Q}_R(M, n) \leq (2 + d \log d(n+1))/n$ for every R in Ω and for all M in the convex hull of \mathcal{D} .

We will prove theorem later in the section.

Because $(n-1)\bar{Q} + w^T w$ is a quadratic over a compact set, the minimum exists but it need not be unique. If it is not unique, just pick one arbitrarily. No matter how this is done, we have

$$\begin{aligned} & (n-1)\bar{Q}((\hat{w}(n-1) + x)^T f, n-1) + (\hat{w}(n-1) + x)^T (\hat{w}(n-1) + x) \\ & - (n-1)\bar{Q}(\hat{w}(n-1)^T f, n-1) - \hat{w}(n-1)^T \hat{w}(n-1) \geq 0 \end{aligned}$$

for all $x \in \mathbb{R}^d$ with $\hat{w}(n-1) + x \in B$. Let $L_{n-1}(x)$ denote the linear part in x of this expression. $L_{n-1}(x)$ can be thought of as the directional derivative of

$(n - 1)\bar{Q} + w^T w$ at its minimum point $(\hat{w}(n - 1))$. If $\hat{w}(n - 1)$ happens to be in the interior B , then $L_{n-1}(x) = 0$. More generally, it follows that $L_{n-1}(x) \geq 0$ if $\hat{w}(n - 1) + x \in B$.

Define Δ , a function of time i and an \mathbb{R}^d -vector x , as

$$\begin{aligned} \Delta(i, x) &= i(\bar{Q}_R(\hat{w}(i - 1)^T f, i) - \bar{Q}_R((\hat{w}(i - 1) + x)^T f, i)) \\ &\quad + \hat{w}(i - 1)^T \hat{w}(i - 1) - (\hat{w}(i - 1) + x)^T (\hat{w}(i - 1) + x). \end{aligned}$$

Using the fact that $\hat{w}(i - 1)$ minimizes $\sum_{j=1}^i (R_j - w^T f)^2 + w^T w$ over B , we have for $\hat{w}(i - 1) + x \in B$,

$$\Delta(i, x) = -x^T G_i x + a_i x^T f(i) - L_{i-1}(x) \leq -x^T G_i x + a_i x^T f(i),$$

where $f(i)$ is the vector of predictions at time i , $G_i = I + \sum_{j=1}^i f(j) f(j)^T$, $a_i = 2(R_i - f(i)^T \hat{w}(i - 1))$ is a real number with $|a_i| \leq 2$ and I is the identity matrix.

Our first claim relates Δ to our total error.

CLAIM A.

$$\begin{aligned} &\bar{Q}_R(S_\emptyset, n) - \bar{Q}_R(\hat{w}(n)^T f, n) \\ &= \frac{1}{n} \left(\hat{w}(n)^T \hat{w}(n) - \hat{w}(0)^T \hat{w}(0) + \sum_{i=1}^n \max_{\hat{w}(i-1)+x \in B} \Delta(i, x) \right). \end{aligned}$$

PROOF.

$$\begin{aligned} \bar{Q}_R(S_\emptyset, n) &= \frac{1}{n} \sum_{i=1}^n \left(i \bar{Q}_R(\hat{w}(i - 1)^T f, i) \right. \\ &\quad \left. - (i - 1) \bar{Q}_R(\hat{w}(i - 1)^T f, i - 1) \right) \\ &= \bar{Q}_R(\hat{w}(n)^T f, n) \\ &\quad + \frac{1}{n} \sum_{i=1}^n i \left(\bar{Q}_R(\hat{w}(i - 1)^T f, i) - \bar{Q}_R(\hat{w}(i)^T f, i) \right) \\ &\quad + \frac{1}{n} \left(\hat{w}(n)^T \hat{w}(n) - \hat{w}(0)^T \hat{w}(0) \right. \\ &\quad \left. - \sum_{i=1}^n \left[\hat{w}(i - 1)^T \hat{w}(i - 1) - \hat{w}(i)^T \hat{w}(i) \right] \right) \\ &= \bar{Q}_R(\hat{w}(n)^T f, n) \\ &\quad + \frac{1}{n} \left(\hat{w}(n)^T \hat{w}(n) - \hat{w}(0)^T \hat{w}(0) \right. \\ &\quad \left. + \sum_{i=1}^n \Delta(i, \hat{w}(i) - \hat{w}(i - 1)) \right). \end{aligned}$$

Notice that $\Delta(i, \cdot)$ attains the maximum at $x = \hat{w}(i) - \hat{w}(i - 1)$. \square

In Claim B, we bound the size of Δ .

CLAIM B.

$$\max_{\hat{w}^{(i-1)} + x \in B} \Delta \leq \max_{x \in \mathbb{R}^d} (-x^T G_i x + a_i x^T f(i)) \leq \frac{|G_i| - |G_{i-1}|}{|G_i|}.$$

PROOF. Note $G_i = G_{i-1} + f(i)f(i)^T$. Because G_i is the identity plus a nonnegative definite matrix, we know that G_i is invertible. Let $|G_i|$ be the determinant of G_i .

$$\begin{aligned} |G_{i-1}| &= |G_i - f(i)f(i)^T|, \\ |G_{i-1}| &= |G_i|(1 - f(i)^T G_i^{-1} f(i)), \\ f(i)^T G_i^{-1} f(i) &= \frac{|G_i| - |G_{i-1}|}{|G_i|}. \end{aligned}$$

But the maximum value of Δ is less than or equal to $a_i f(i)^T G_i^{-1} f(i) a_i / 4$. (Maximizing the quadratic $-x^T A x + b^T x$, where A is a positive definite matrix, yields $b^T A^{-1} b / 4$.) Recalling that $|a_i| \leq 2$, we get the desired result. \square

CLAIM C.

$$\sum_{i=1}^n \max_{\hat{w}^{(i-1)} + x \in B} \Delta(i, x) \leq d \log d(n + 1).$$

PROOF. Using Claim B and the equation $1 - x \leq -\log(x)$, we get

$$\sum_{i=1}^n \max \Delta(i, x) \leq \log(|G_n|).$$

$|G_n| \leq L^d$, where L is the maximum eigenvalue of G_n . The elements of G_n are bounded by $n + 1$. Thus, L is bounded by $d(n + 1)$. \square

PROOF OF THEOREM 1. For all M in the convex hull of \mathcal{D} ,

$$\begin{aligned} &\bar{Q}_R(M, n) - \bar{Q}_R(S_{\mathcal{D}}, n) \\ &\geq \bar{Q}_R(\hat{w}(n)^T f, n) - \bar{Q}_R(S_{\mathcal{D}}, n) - \frac{1}{n} \\ &\geq \frac{1}{n} \left(- \sum_{i=1}^n \max \Delta(i, x) + \hat{w}(0)^T \hat{w}(0) - \hat{w}(n)^T \hat{w}(n) - 1 \right) \\ &\geq \frac{1}{n} (-d \log d(n + 1) - 2), \end{aligned}$$

where the first inequality follows from $\bar{Q}_R(\hat{w}(n)^T f, n) \leq \bar{Q}_R(M, n) + 1/n$ for all M in \mathcal{D} . The second line is by Claim A. For $w \in B$, $w^T w \leq 1$. This plus Claim C yields the last inequality. Thus

$$d_R(M, S_{\mathcal{D}}) \geq \liminf_{n \rightarrow \infty} \frac{-d \log d(n+1) - 2}{n} = 0. \quad \square$$

Since f_i is in the convex hull of \mathcal{D} , we see that $S_{\mathcal{D}}$ is as strong as f_i . So, of the predictors, $\{f_1, f_2, \dots, f_d, S_{\mathcal{D}}\}$, the strongest is $S_{\mathcal{D}}$.

The bound in Theorem 1 holds for all R , so in particular it holds for the worst possible R —called the worst case. Thus it deals with prediction in the worst case.

THEOREM 2. *For any countable set \mathcal{D} of predictors, there exists a predictor $T_{\mathcal{D}}$ which is as strong as every member of \mathcal{D} .*

PROOF. Define T_i for $i = 1, 2, 3, \dots$ as

$$T_i(n) = \begin{cases} f_i(n), & \text{if } i \geq n, \\ S_{\mathcal{D}_i}(n), & \text{if } i < n, \end{cases}$$

where $\mathcal{D}_i = \{T_{i+1}, f_i\}$. Notice that for all times n , all T_i are well-defined. By Theorem 1, T_i is as strong as f_i and as strong as T_{i+1} . By transitivity, T_1 is as strong as f_i . In other words, T_1 is as strong as T_2 , which is as strong as T_3, \dots , which is as strong as T_i , which is as strong as f_i , so T_1 is as strong as f_i . Letting $T_{\mathcal{D}} = T_1$ proves our result. \square

EXAMPLE. Theorem 2 is only an asymptotic result. In other words, a set \mathcal{D} can be constructed such that at all finite times n , there exist a perfect $f_i \in \mathcal{D}$ with $\bar{Q}(f_i, n) = 0$, but $\bar{Q}(T_{\mathcal{D}}, n) \cong \frac{1}{2}$. Theorem 2 forces $\liminf_{n \rightarrow \infty} \bar{Q}(f_i, n) \geq \frac{1}{2}$ for all f_i in \mathcal{D} . So, in this case,

$$0 = \liminf_{n \rightarrow \infty} \bar{Q}(f_i, n) < \liminf_{n \rightarrow \infty} \bar{Q}(T_{\mathcal{D}}, n) \leq \inf_i \liminf_{n \rightarrow \infty} \bar{Q}(f_i, n).$$

Acknowledgments. I am grateful for the help of Grace Yang, Michael Cohen, Rick Vohra and the referees. But the most thanks go to Peyton Young who provided encouragement, support and a general guiding hand.

REFERENCES

BRIER, G. (1950). Verification of forecasts expressed in terms of probability. *Monthly Weather Review* **78** 1–3.
 COVER, T. M. (1991). Universal portfolios. *Math. Finance* **1** 1–12.
 DAWID, A. P. (1984). Statistical theory: The prequential approach (with discussion). *J. Roy. Statist. Soc. Ser. A* **147** 278–292.

- DAWID, A. P. (1986). Probability forecasting. In *Encyclopedia of Statistical Sciences* 7 210–218. Wiley, New York.
- GUPTA, S. and WILTON, P. C. (1987). Combination of forecasts: An extension. *Management Sci.* 3 356–372.
- NEWBOLD, P. and GRANGER, C. W. J. (1974). Experience with forecasting univariate time series and the combination of forecasts. *J. Roy. Statist. Soc. Ser. A* 137 131–146.
- RISSANEN, J. (1986). Stochastic complexity and modeling. *Ann. Statist.* 14 1080–1100.

UNIVERSITY OF CHICAGO
GRADUATE SCHOOL OF BUSINESS
CHICAGO, ILLINOIS 60637