

Efficient and Truthful Forecasting Competitions

Rafael Frongillo, Robert Gomez, Anish Thilagar, Bo Waggoner

January 12, 2021

Abstract

Witkowski et al. [2018] studied competitions among forecasters who predict a set of common events. When such a competition picks a single winner, incentives can be distorted, so that paper introduced the ELF mechanism which incentivizes truthful forecasts. In this paper, we give a tighter analysis of the ELF mechanism, then give a mechanism that improves upon it exponentially, in the sense of also being truthful but requiring only logarithmically many events in order to select the best forecaster with high probability. These results have potential applications to machine-learning competitions such as Kaggle, where winners are currently selected in non-truthful ways.

1 Introduction

2 Model and Definitions

Our model will consist of n forecasters s.t. $n \geq 2$ where $i \in [1, n]$ denotes the i th forecaster, m independent events where $k \in [1, m]$ denotes the k th event. Each event $k \in [1, m]$ will be associated with a random variable X_k . $X_k = 1$ will denote that event k occurred, and, likewise, $X_k = 0$ will denote event k did not occur. $\forall k \in [m] \exists \theta_k$ s.t. θ_k denotes the true, unknown probability of event k occurring;

$\theta_k = P(X_k = 1) = 1 - P(X_k = 0)$. We will denote each outcome drawn as $x_i \sim X_i$. For each forecaster i , we denote their true belief that event k will occur as $\hat{\theta}_{i,k} \in (0, 1)$. We then denote $y_{i,k} \in [0, 1]$ as the actual report forecaster i makes for event k . Note that $y_{i,k}$ need not equal $\hat{\theta}_{i,k}$.

After every event $k \in [1, m]$ has occurred, and thus each forecaster has reported $y_{i,k} \forall i \forall k$, a mechanism will select exactly one forecaster from the set of $[n]$ as the 'best forecaster'.

Definition 1. A **selection mechanism** $M(y, x)$ is a map from the forecasts $y = (y_1, \dots, y_n)$ and outcomes of the events $x = (x_1, \dots, x_k)$ to an index $w \in [1, n]$.

We say that forecaster w wins the competition if $M(y, x) = w$. We call a forecaster selection mechanism truthful iff, for all forecasters, their probability of winning is maximized when they reports their beliefs as their forecasts.

Definition 2. M is **truthful** iff $\forall i \in [n], \forall \hat{\theta}_i, \forall y'_i \neq \hat{\theta}_i$, and $\forall y_j \text{ s.t. } j \neq i$,

$$P(M((y_1, \dots, \hat{\theta}_i, \dots, y_n), x) = i) \geq P(M((y_1, \dots, \hat{y}_i, \dots, y_n), x) = i)$$

M is **strictly truthful** if the strict inequality inequality holds.

Additionally, we would like our mechanism to choose the forecaster with the best predictions. We define the accuracy of a single forecaster i as 1 minus the average squared loss of their forecasts compared to the true probabilities.

Definition 3. A forecaster's **accuracy** is $a_i = 1 - \frac{1}{m} \sum_{k=1}^m (y_{i,k} - \theta_k)^2$.

An ideal mechanism would choose forecasters with larger accuracies with higher probabilities. Let $i^* = \arg \max_i a_i$ be the highest accuracy forecaster, and let $a^* = a_{i^*}$.

Definition 4. The **accuracy gap** $\epsilon = a^* - \max_{i \neq i^*} a_i$ is the difference in accuracy of the first and second best forecasters.

For a given mechanism M and n forecasters with accuracy gap ϵ , we let m_δ be the minimum number of events such that M necessarily chooses i^* with probability at least $1 - \delta$. Specifically, for any $m > m_\delta(M, n, \epsilon)$, $P(M(y, x) \neq i^*) \leq \delta$. m_δ gives us a measure of how many events are needed for a mechanism to pick the best forecaster. We will use it as a metric to compare mechanisms, where we say M is more efficient than M' if $m_\delta(M, n, \epsilon) < m_\delta(M', n, \epsilon)$, since M needs less events to consistently pick the best forecaster.

Similar to work in Witkowski et al. [2018], we will focus on using the proper scoring rule, the quadratic scoring rule in our forecaster selection mechanism.

Definition 5. A **scoring rule**, R is proper if $\forall p, y \in [0, 1], E_{x \sim p}[R(p, x)] \geq E_{x \sim p}[R(y, x)]$

Definition 6. The *quadratic scoring rule*, is $R_q = 1 - (y - x)^2$, where x denotes the observed outcome, and y denotes the report.

2.1 Useful facts

Bo: TODO name this subsection better.

The following facts, easy to derive and used by Witkowski et al. [2018], will be useful to us as well.

Lemma 7. *The expected score, from the ground truth point of view, of forecaster i on event k is $1 - \theta_k(1 - \theta_k) - (y_{i,k} - \theta_k)^2$.*

Proof. By definition, the expected score is

$$\begin{aligned} & \theta_k [1 - (1 - y_{i,k})^2] + (1 - \theta_k) [1 - y_{i,k}^2] \\ &= \theta_k [2y_{i,k} - y_{i,k}^2] + (1 - \theta_k) [1 - y_{i,k}^2] \\ &= 1 - \theta_k + 2\theta_k y_{i,k} - y_{i,k}^2 \\ &= 1 - \theta_k + \theta_k^2 - (y_{i,k} - \theta_k)^2. \end{aligned}$$

□

Averaging over events, we get that i 's score on average over the m events, in expectation, is exactly her accuracy minus some penalty term that only depends on how extreme the ground truth events are.

Corollary 8. *The expectation of forecaster i 's average score on all events, from the ground truth point of view, is $a_i - \frac{1}{m} \sum_{k=1}^m \theta_k(1 - \theta_k)$, where a_i is i 's accuracy.*

We also immediately obtain:

Corollary 9. *The expected difference in average scores of forecaster i and j , from a ground truth point of view, is $a_i - a_j$.*

3 Simple average mechanism

A straightforward selection mechanism often used in practice is what we call the *simple average mechanism* M_s . For this mechanism, we assign each forecaster a score $f_{i,k} = R_q(y_{i,k}, x_k)$ for each event. Then, we assign their final score by summing

these over all events, $F_i = \sum_{j=1}^m f_{i,j}$. Finally, we choose the forecaster with the highest cumulative score $M_s(y, x) = \arg \max_i F_i$.

3.1 Truthfulness

As observed by Witkowski et al. [2018] and discussed in depth by Lichtendahl and Winkler [2007], this mechanism is generally not truthful. We can consider a simple example with three forecasters and one event, whose true probability is 0.5. Alice predicts 0.5, but Bob predicts 1 and Charlie predicts 0. Observe that Alice *cannot win*, despite being the best forecaster by far: either the event occurs (Bob wins) or it doesn't (Charlie wins). She can only win by predicting either 0 or 1, raising her probability of winning from 0 to 0.25, assuming a random tiebreaking rule.

3.2 Efficiency

Although the simple average mechanism is not truthful, it is worth studying its efficiency because, intuitively, it seems unlikely to do better using a mechanism that relies on the quadratic scores of the participants. So as a benchmark for truthful quadratic-score-based mechanisms, let us consider the efficiency of M_s if forecasters were to report their true beliefs $\hat{\theta}_i$.

For a single event, the most extreme forecasters will always get selected, as in the example above. However, for multiple events, the averaging starts to favor more accurate forecasters. In fact, one can show that with n forecasters and an accuracy gap of ϵ , only

$$m = \frac{2}{\epsilon^2} \ln \left(\frac{n}{\delta} \right)$$

events suffice to select the best forecaster with probability $1 - \delta$.

Sketch (similar to the Witkowski et al. [2018] analysis of ELF): Let i be the best forecaster. She beats some competitor j if she has a higher average score in the competition. The expected difference in average scores is the difference in accuracies (Corollary 9), which is $a_i - a_j \geq \epsilon$. On any given event, the difference in scores is an independent random variable bounded in $[-1, 1]$. A standard Hoeffding bound shows that the average of $m = \frac{2}{\epsilon^2} \ln \frac{n}{\delta}$ such variables is at least ϵ below its mean

with probability $1 - \frac{\delta}{n}$. Complete the argument by taking a union bound over all $n - 1 \leq n$ opponents j .

This shows that it is possible to select the best forecaster using a number of events only logarithmic in n , the number of forecasters – if truthfulness is not required. We next consider the number of events used by the ELF mechanism.

4 ELF

The ELF mechanism M_{ELF} [Witkowski et al., 2018], improves upon M_s , by incentivizing forecasters to be truthful. For each event k , we hold a lottery to choose a winner. Forecaster i is chosen with probability

$$f_{i,k} = \frac{1}{n} + \frac{1}{n} \left(R_q(y_{i,k}, x_k) - \frac{\sum_{j \neq i} R_q(y_{j,k}, x_k)}{n-1} \right)$$

This score is the additively normalized quadratic score among all forecasters. A forecaster increasing their own quadratic score increases the chances they win an event, while decreasing the chances others win that event. Let w_k be the winner of event k . Then we choose the forecaster $M_{\text{ELF}}(y, x) = \arg \max_i \sum_{k=1}^m \mathbb{1}(w_k = i)$ as the winner of the competition. Essentially, we give the winner of each event lottery a point, and then pick the forecaster who got the most points as the overall winner.

In Theorem 6 of Witkowski et al. [2018], it is proven that M_{ELF} is strictly truthful. In Theorem 8 of the same paper, they also upper bound m_δ by showing

$$m_\delta(M_{\text{ELF}}, n, \epsilon) \leq \frac{2(n-1)^2}{\epsilon^2} \ln \left(\frac{4(n-1)}{\delta} \right)$$

For a fixed forecaster accuracy gap and desired mechanism accuracy, the number of events ELF needs to consistently choose the right forecaster is proportional to $\frac{1}{\epsilon^2} n^2 \ln n$. However, this upper bound is not tight.

Let $F_{i,k}$ be the indicator function for forecaster i winning the lottery for event k . Then, $F_i = \sum_k F_{i,k}$ is the random variable that is the number of events forecaster i wins. The proof of Theorem 8 Witkowski et al. [2018] utilized the Hoeffding bound applied to F_i , which assumes they have high variance. However, this is not the case, and utilizing that fact gives an improvement in the lower bound.

Theorem 10 (Bernstein's Inequality). *Given random variables X_i for $1 \leq i \leq n$ such that $0 \leq X_i \leq 1$ almost surely, let $S = \sum_i X_i$. Then,*

$$P(|S - \mathbb{E}[S]| < t) < 2e^{\frac{-t^2}{2(\text{Var}(S) + \frac{t}{3})}}$$

Theorem 11. $m_\delta(M_{\text{ELF}}, n, \epsilon) \leq \frac{5(n-1)}{\epsilon^2} \ln \left(\frac{4(n-1)}{\delta} \right)$.

Proof: Note that for any single event k , every forecaster wins with probability at best $\frac{2}{n}$, and at worst $\frac{1}{n} - \frac{1}{n(n-1)}$ (achieved when $y_i = 1$, $y_{j \neq i} = 0$, and $x_i = 1$). Since $F_{i,k}$ is just a binary random variable, it will have highest variance when its expectation is as close to $\frac{1}{2}$ as possible. For nontrivially small n , this will mean its expectation is $\frac{2}{n}$ (or equivalently $P(F_{i,k} = 1) = \frac{2}{n}$), so its variance will be at most $\frac{2(n-2)}{n^2}$. So we have $\text{Var}(F_{i,k}) \leq \frac{2(n-2)}{n^2} < \frac{2}{n}$.

Since $F_i = \sum_k F_{i,k}$, we have $\text{Var}(F_i) = \sum_k \text{Var}(F_{i,k}) < \sum_k \frac{2}{n} = \frac{2m}{n}$. As shown in Witkowski et al. [2018], for $j \neq i$ we have

$$\mathbb{E}[F_i] - \mathbb{E}[F_j] \geq \frac{m\epsilon}{n-1}$$

Therefore, if $F_j \geq F_i$, (forecaster j beats forecaster i), then either forecaster j overperformed their expectation by at least $\frac{m\epsilon}{2(n-1)}$, or forecaster i underperformed their expectation by $\frac{m\epsilon}{n-1}$. Specifically, either $\mathbb{E}[F_i] - F_i \geq \frac{m\epsilon}{2(n-1)}$ or $F_j - \mathbb{E}[F_j] \geq \frac{m\epsilon}{2(n-1)}$

Applying Theorem 10, we have

$$\begin{aligned} P\left(|F_i - \mathbb{E}[F_i]| < \frac{m\epsilon}{2(n-1)}\right) &< 2e^{\frac{-\left(\frac{m\epsilon}{2(n-1)}\right)^2}{2\left(\frac{2m}{n} + \frac{m\epsilon}{6(n-1)}\right)}} \\ &< 2e^{\frac{-\left(\frac{m\epsilon}{2(n-1)}\right)^2}{2\left(\frac{m}{n-1}\left(2 + \frac{\epsilon}{6}\right)\right)}} \\ &= 2e^{\frac{-\frac{m\epsilon^2}{2(n-1)}}{\left(2 + \frac{\epsilon}{6}\right)}} \\ &< 2e^{-\frac{m\epsilon^2}{5(n-1)}} \end{aligned} \tag{1}$$

And similarly for F_j . Putting both of those cases together, we have by the union

bound

$$\begin{aligned}
P(F_j \geq F_i) &= P\left(\left(F_j - \mathbb{E}[F_j] < \frac{m\epsilon^2}{2(n-1)}\right) \cup \left(\mathbb{E}[F_i] - F_i < \frac{m\epsilon^2}{2(n-1)}\right)\right) \\
&\leq P\left(F_j - \mathbb{E}[F_j] < \frac{m\epsilon^2}{2(n-1)}\right) + P\left(\mathbb{E}[F_i] - F_i < \frac{m\epsilon^2}{2(n-1)}\right) \\
&\leq 4 e^{-\frac{m\epsilon^2}{5(n-1)}}
\end{aligned}$$

Using the union bound over all j , we have

$$\begin{aligned}
P(M_l(y_1, \dots, y_n, x) = i) &= 1 - \sum_{j \neq i} P(M_l(y_1, \dots, y_n, x) = j) \\
&\geq 1 - \sum_{j \neq i} P(F_j \geq F_i) \\
&\geq 1 - 4(n-1)e^{-\frac{m\epsilon^2}{5(n-1)}}
\end{aligned}$$

Assigning this probability to $1 - \delta$, we see that for a fixed n, ϵ , ELF will choose the correct forecaster with probability at least $1 - \delta$ if m is large enough for the above inequality to hold. Solving it for m , we get the minimum such satisfying value

$$m_\delta \leq \frac{5(n-1)}{\epsilon^2} \ln \left(\frac{4(n-1)}{\delta} \right)$$

□

Theorem 12. For $\delta < \frac{1}{2}$, $m_\delta(M_{\text{ELF}}, n, \epsilon) > n \log n$

Proof: Pick any $\delta < \frac{1}{2}$, $0 < c \leq 1$ and consider the case where $m = cn \log n$ and $\theta_k = 1$ for every event. Let $y_{1,k} = 1$ and $y_{i,k} = 0$ for every $i > 1$. Therefore, forecaster 1 always forecasts correctly and has a quadratic score $R_q(1, 1) = 1$ for each event, while everyone else is wrong and has a quadratic score of $R_q(0, 1) = 0$ for each event. Then, $a_1 = 1$, and $a_i = 0$ for $i > 1$, so we have the worst possible accuracy gap $\epsilon = 1$.

For every event, forecaster 1's score will be

$$\begin{aligned}
f_1 &= \frac{1}{n} \left(1 + R_q(y_1, 1) - \frac{1}{n-1} \sum_{j \neq i} R_q(y_j, x) \right) \\
&= \frac{1}{n} \left(1 + 1 - \frac{1}{n-1} \sum_{j \neq i} 0 \right) \\
&= \frac{2}{n}
\end{aligned}$$

Since all the other forecasters are symmetric, their f_i must all be the same, and since they sum to 1, we must have $f_i = \frac{n-2}{n(n-1)} < \frac{1}{n}$.

We can reframe the way winners are chosen for each event. Instead of holding a lottery, we first flip a coin that gives forecaster 1 a win with probability $\frac{2}{n}$. If it does not give them the win, we run a normal lottery for the remaining forecasters, where they each have probability $\frac{1}{n-1}$ of winning. This gives them an overall probability of $\frac{n-2}{n(n-1)}$ to win each event, so it is the same as the original lottery.

Since all the events are uniform, the expected number of lotteries forecaster 1 will win is $mf_1 = \frac{2m}{n} = 2c \log n$. Specifically, their wins will follow the Bernoulli distribution $B(m, \frac{2}{n})$. With some probability $C_1 = P(B(m, \frac{2}{n}) < 2c \log n) \approx \frac{1}{2}$, forecaster 1 will win less events than expected.

The distribution of wins for the remaining forecasters will follow a uniform multinomial distribution. Specifically, we can model it with Balls and Bins. Let M be the maximum number of wins of any of the remaining forecasters.

Lemma 13. $P(M > 2c \log n) = 1 - o(1)$ for any $0 < c \leq 1$.

Proof: From Theorem 1 of Raab and Steger [1998], we know that for $m = cn \log n$, $P(M > k_\alpha) = 1 - o(1)$, where $k_\alpha = (d_c - 1 + \alpha) \log n$, and $0 < 1 < \alpha$, such that d_c is the solution of $1 + x(\log c - \log x + 1) - c = 0$ such that $x > c$.

If $c = 1$, this reduces to $d_1(1 - \log d_1) = 0$, so $d_1 = e$. Since this inequality holds for any α , we let $\alpha = 3 - e$. Plugging in both values gives us $k_\alpha = 2 \log n$, so $P(M > 2 \log n) = 1 - o(1)$.

We can then rewrite the equation that determines d_c as $1 + \frac{1}{x} = \frac{c}{x} - \log \frac{c}{x}$. $\frac{d}{dx}$ of the left hand side is $-\frac{1}{x^2} < 0$ while the right hand side is $\frac{x-c}{x^2} > 0$, since we are only considering the solution where $x > c$. For any $c < 1$, $x = ce$ will not satisfy the equation, since the right hand side remains constant while the left is increasing from our d_1 solution. In order for the left side to decrease and the right side to increase, x must increase. Therefore, $d_c > ce$ for $0 < c < 1$.

Letting $\alpha = \max(c(2 - e) + 1, 0)$, we have

$$k_\alpha = (d_c - 1 + \alpha) \log n > (ce - 1 + c(2 - e) + 1) \log n = 2c \log n.$$

We can again apply Theorem 1 of Raab and Steger [1998] to get

$$P(M > 2c \log n) > P(M > k_\alpha) = 1 - o(1)$$

for any $c < 1$. Therefore, for any $c \leq 1$, $P(M > 2c \log n) = 1 - o(1)$. \square

Using the lemma, for any $C_2 < 1$ and sufficiently large n , $P(M > 2c \log n) > C_2$. We choose some $C_2 > 2\delta$. Now, the probability that forecaster 1 is not chosen is at least

$$P\left(B\left(m, \frac{2}{n}\right) < 2c \log n\right) \times P(M > 2c \log n) > \frac{1}{2} \times 2\delta = \delta$$

Therefore, it is not true that for any n, ϵ , $P(M_{\text{ELF}}(y, x) \neq i^*) \leq \delta$ when $m = cn \log n$. Therefore, $m_\delta(M_{\text{ELF}}, n, \epsilon) \neq cn \log n$ for any $c < 1$.

Therefore, $m_\delta(M_{\text{ELF}}, n, \epsilon) > n \log n$. □

4.1 Other point-per-round mechanisms

We can extend Theorem 12 to apply to any normal, anonymous mechanism that assigns a single point per event to 1 forecaster, and then picks the forecaster with the highest score. We can treat such a forecaster selection mechanism as a budget-balanced wagering mechanism for each event, where each forecaster has a shared, fixed wager of $\frac{1}{n}$. Since we want the mechanism to be symmetric, we can assume it satisfies anonymity as described in [wagering mechanism paper].

(TODO: justify normality).

Then, by Lemma 4 in [wagering mechanism paper], we have that forecaster i 's payout for an event will be

$$\frac{1}{n} + f(y_i, x) - \frac{1}{n-1} \sum_{j \neq i} f(y_j, x)$$

where f is a proper scoring rule whose values are contained in an interval of size $\frac{1}{n}$.

References

- Kenneth C. Lichtendahl, Jr. and Robert L. Winkler. Probability elicitation, scoring rules, and competition among forecasters. *Management Science*, 53(11):1745–1755, 2007.
- Martin Raab and Angelika Steger. “balls into bins” - a simple and tight analysis. In *Proceedings of the Second International Workshop on Randomization and Approximation Techniques in Computer Science*, RANDOM 1998, pages 159–170. Springer-Verlag, 1998. ISBN 3-540-65142-X.

Jens Witkowski, Rupert Freeman, Jennifer Wortman Vaughan, David M. Pennock, and Andreas Krause. Incentive-compatible forecasting competitions. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence*, AAAI 2018, 2018.