

On the multibin logarithmic score used in the FluSight competitions

Johannes Bracher^{a,1}

The FluSight challenges (1) represent an outstanding collaborative effort and have “pioneered infectious disease forecasting in a formal way” (ref. 2, p. 2803). However, I wish to initiate a discussion about the employed evaluation measure.

The competitions feature discrete or discretized targets related to the US influenza season. E.g., for the peak timing Y , a forecast distribution F consists of probabilities p_1, \dots, p_T for the $T = 33$ wk of the season. Such forecasts can be evaluated using the log score (3, 4)

$$\log S(F, y_{\text{obs}}) = \log(p_{y_{\text{obs}}}),$$

where y_{obs} is the observed value. This score is strictly proper; i.e., its expectation is uniquely maximized by the true distribution of Y . In the FluSight competitions the logS is applied in a multibin version,

$$\text{MBlogS}(F, y_{\text{obs}}) = \log\left(\sum_{i=-d}^d p_{y_{\text{obs}}+i}\right),$$

to measure “accuracy of practical significance” (ref. 1, p. 3153). Depending on the target, d is either 1 or 5. Following the competitions, this score has become widely used (5–10), even though as also mentioned in ref. 1, it is improper. This may be problematic as improper scores incentivize dishonest forecasts. Assume $T > 2d$ and

$$p_1 = \dots = p_d = p_{T-d+1} = \dots = p_T = 0, \quad [1]$$

i.e., probability 0 for the $2d$ extreme categories. Now define a “blurred” distribution \tilde{F} with

$$\tilde{p}_t = \frac{\sum_{i=-d}^d p_{t+i}}{2d+1}, t = 1, \dots, T, \quad [2]$$

where $p_t = 0$ for $t < 1$ and $t > T$ and Eq. 1 ensures $\sum_{t=1}^T \tilde{p}_t = 1$. This implies

$$\text{MBlogS}(F, y_{\text{obs}}) = \log S(\tilde{F}, y_{\text{obs}}) + \log(2d+1);$$

i.e., the MBlogS is essentially the logS applied to a blurred version of F . To optimize the expected MBlogS under their true belief F , forecasters should therefore not report F , but a sharper forecast G so that the blurred version \tilde{G} (with $\tilde{p}_{G,1}, \dots, \tilde{p}_{G,T}$ derived from $p_{G,1}, \dots, p_{G,T}$ as in Eq. 2) is close or equal to F . This follows from the propriety of the logS. An optimal G is found by maximizing $\sum_{t=1}^T p_t \cdot \log(\tilde{p}_{G,t})$ with respect to $p_{G,1}, \dots, p_{G,T}$.

This optimal G can differ considerably from the original F , as Fig. 1 shows for forecasts of the 2016 to 2017 national-level peak timing by the Los Alamos National Laboratory (LANL) team (9) (downloaded from <https://github.com/FluSightNetwork/cdc-flusight-ensemble/>). The optimized G s (with $d = 1$) often have their mode shifted by 1 wk and tend to be multimodal, even for unimodal F . Averaged over the 2016 to 2017 season they yield improved MBlogS for the peak timing (−0.434 vs. −0.484). This illustrates that the MBlogS may be gamed, even though we strongly doubt participants have tried to. The logS, like any other proper score, could avoid such pitfalls.

Acknowledgments

I thank T. Gneiting for helpful discussions and the FluSight Collaboration for making its forecasts publicly available.

^aEpidemiology, Biostatistics and Prevention Institute, University of Zurich, 8001 Zurich, Switzerland

Author contributions: J.B. designed research, performed research, analyzed data, and wrote the paper.

The author declares no conflict of interest.

Published under the [PNAS license](#).

¹Email: johannes.bracher@uzh.ch.

First published September 26, 2019.

