

¹ Transformation of forecasts for evaluating predictive performance in
² an Scoring epidemiological context forecasts on transformed scales

³ Nikos I. Bosse^{1,2,3,*}, Sam Abbott^{1,2}, Anne Cori⁴,
Edwin van Leeuwen^{1,5}, Johannes Bracher^{6,7,†}, Sebastian Funk^{1,2,3,†}

⁴ June 26, 2023

⁵ **Abstract**

Forecast evaluation plays an essential role in the development cycle is essential for the development of predictive epidemic models and can inform their use for public health decision-making. Common scores to evaluate epidemiological forecasts are the Continuous Ranked Probability Score (CRPS) and the Weighted Interval Score (WIS), which are both can be seen as measures of the absolute distance between the forecast distribution and the observation. They are commonly applied However, applying these scores directly to predicted and observed incidence counts , but it can be questioned whether this is the optimal procedure for comparing models given may not be the most appropriate due to the exponential nature of epidemic processes and the several orders of magnitude that observed values can span over varying magnitudes of observed values across space and time. In this paper, we argue that transforming counts before applying scores such as the CRPS or WIS can effectively mitigate these difficulties and yield epidemiologically meaningful and easily interpretable results. We motivate the procedure threefold using the Using the CRPS on log-transformed counts-values as an example, we list three attractive properties: Firstly, it can be interpreted as a probabilistic version of a relative error. Secondly, it reflects how well models predicted the time-varying epidemic growth rate. And lastly, using arguments on variance-stabilizing transformations, it can be shown that under the assumption of a quadratic mean-variance relationship, the logarithmic transformation leads to expected CRPS values which are independent of the order of magnitude of the predicted quantity. Applying the log transformation a transformation of $\log(x + 1)$ to data and forecasts from the European COVID-19 Forecast Hub, we find that it changes model rankings regardless of stratification by forecast date, location or target types. Situations in which models missed the beginning of upward swings are more strongly emphasized emphasised while failing to predict a downturn following a peak is less severely penalized penalised when scoring transformed forecasts as opposed to untransformed ones. We conclude that appropriate transformations, of which the natural logarithm is only one particularly attractive option, should be considered when assessing the performance of different models in the context of infectious disease incidence.

³⁰ * Correspondence to nikos.bosse@lshtm.ac.uk, † Contributed equally

¹Department of Infectious Disease Epidemiology, London School of Hygiene & Tropical Medicine, London, United Kingdom

²Centre for the Mathematical Modelling of Infectious Diseases, London, United Kingdom

³NIHR Health Protection Research Unit in Modelling & Health Economics

⁴MRC Centre for Outbreak Analysis and Modelling, Department of Infectious Disease Epidemiology, School of Public Health, Imperial College London, London, United Kingdom

⁵Modelling & Economics Unit and NIHR Health Protection Research Unit in Modelling & Health Economics, UK Health Security Agency, London, United Kingdom

⁶Chair of Statistical Methods and Econometrics, Karlsruhe Institute of Technology, Karlsruhe, Germany

⁷Computational Statistics Group, Heidelberg Institute for Theoretical Studies, Heidelberg, Germany

31 1 Introduction

32 Probabilistic forecasts (Held et al., 2017) play an important role in decision-making in epidemiology and
 33 public health (Reich et al., 2022), as well as other areas as diverse as economics (Timmermann, 2018) or
 34 meteorology (Gneiting and Raftery, 2005). Forecasts based on epidemiological modelling in particular ~~has~~
 35 ~~have~~ received widespread attention during the COVID-19 pandemic. Evaluations of forecasts can provide
 36 feedback for researchers to improve their models and train ensembles. They moreover help decision-makers
 37 distinguish good from bad predictions and choose forecasters and models that are best suited to inform
 38 future decisions.

39 Probabilistic forecasts are usually evaluated using so-called ~~(strictly)~~ proper scoring rules (Gneiting and
 40 Raftery, 2007), which return a numerical score as a function of the forecast and the observed data. Proper
 41 scoring rules are constructed such that ~~forecasters they encourage honest forecasting and cannot be ‘gamed’~~
 42 or ‘cheated’. ~~Assuming that the forecaster’s actual best judgement corresponds to a predictive distribution~~
 43 ~~F, a proper score is constructed such that if F was the data-generating process, no other distribution G~~
 44 ~~would yield a better expected score. A scoring rule is called *strictly* proper if there is no other distribution~~
 45 ~~that under F achieves the same expected score as F, meaning that any deviation from F leads to a worsening~~
 46 ~~of expected scores. Forecasters (anyone or anything that issues a forecast) are thus incentivised to report~~
 47 ~~their true belief F about the future. Examples of Common proper scoring rules that have been used to assess~~
 48 ~~epidemiological forecasts are the Continuous Ranked Probability Score (CRPS, Gneiting and Raftery, 2007)~~
 49 ~~or its discrete equivalent, the Ranked Probability Score (RPS, Funk et al., 2019), and are the Weighted~~
 50 ~~Interval Score (Bracher et al., 2021a). logarithmic or log score (Good, 1952) and the continuous ranked~~
 51 ~~probability score (CRPS, Gneiting and Raftery, 2007). The log score is the predictive log density or probability~~
 52 ~~mass evaluated at the observed value. It is supported by the likelihood principle (Winkler, 1996) and has~~
 53 ~~many desirable theoretical properties; however, the particularly severe penalties it assigns to occasional~~
 54 ~~misguided forecasts make it little robust (Bracher et al., 2021a). Moreover, it is not easily applied to~~
 55 ~~forecasts reported as samples or quantiles, as used in many recent disease forecasting efforts. It is nonetheless~~
 56 ~~occasionally used in epidemiology (see e.g., Held et al. 2017; Johansson et al. 2019), but in recent years the~~
 57 ~~CRPS and the weighted interval score (WIS, Bracher et al., 2021a) have become increasingly popular.~~

58 The CRPS measures the distance of the predictive distribution to the observed data as

$$59 \quad \text{CRPS}(F, y) = \int_{-\infty}^{\infty} (F(x) - \mathbf{1}(x \geq y))^2 dx,$$

$$60 \quad \text{CRPS}(F, y) = \int_{-\infty}^{\infty} (F(x) - \mathbf{1}(x \geq y))^2 dx, \quad (1)$$

62 where y is the true observed value and F is the cumulative distribution function (CDF) of the predictive
 63 distribution, and $\mathbf{1}()$ is the indicator function. The CRPS can be understood as a generalisation of the
 64 absolute error to predictive distributions, and interpreted on the natural scale of the data. The WIS is an
 65 approximation of the CRPS for predictive distributions represented by a set of predictive quantiles and is
 66 currently used to assess forecasts in the so-called COVID-19 Forecast Hubs in the US (Cramer et al., 2020,
 67 2021), Europe (Sherratt et al., 2022) and Germany and Poland (Bracher et al., 2021b, 2022), as well as the
 68 US Influenza Forecasting Hub FluSight project on influenza forecasting (CDC, 2022). The WIS is defined as

$$69 \quad \text{WIS}(F, y) = \frac{1}{K} \times \sum_{k=1}^K 2 \times [\mathbf{1}(y \leq q_{\tau_k}) - \tau_k] \times (q_{\tau_k} - y),$$

$$70 \quad \text{WIS}(F, y) = \frac{1}{K} \times \sum_{k=1}^K 2 \times [\mathbf{1}(y \leq q_{\tau_k}) - \tau_k] \times (q_{\tau_k} - y), \quad (2)$$

72 where q_τ is the τ quantile of the forecast F , y is the observed outcome, and K is the number of (roughly
 73 equally spaced) predictive quantiles provided and $\mathbf{1}$ is the indicator function. The WIS can be decomposed
 74 into three components, dispersion, overprediction, underprediction, underprediction and overprediction,
 75 which reflect the width spread of the forecast and whether it was centred above or below the observed
 76 value. We show an alternative definition based on central prediction intervals in [Section Supplement A.1](#)
 77 which illustrates this decomposition.

78 The dynamics of infectious processes are often described by the complementary concepts of the reproduction
 79 number R (Gostic et al., 2020) and growth rate r (Wallinga and Lipsitch, 2007), where R describes the
 80 strength and r the speed of epidemic growth (Dushoff and Park, 2021). In the absence of changes in
 81 immunity, behaviour or other factors that may affect the intensity of transmission, the reproduction number
 82 would be expected to remain approximately constant. In that case, the number of new infections in the
 83 population grows exponentially in time. This behaviour was observed, for example, early in the COVID-19
 84 pandemic in many countries (Pellis et al., 2021).

85 If case numbers are evolving based on an exponential process and the notion of absolute distance encoded
 86 by the CRPS and WIS provides a straightforward interpretation, but may not always be the most useful
 87 perspective in the context of infectious disease spread. Especially in their early phase, outbreaks are best
 88 conceived as exponential processes, characterized by potentially time varying reproduction numbers R_t
 89 (Gostic et al., 2020) or epidemic growth rates r_t (Dushoff and Park, 2021). If the true modelling task re-
 90 volves around estimating and forecasting the reproduction number or the corresponding growth rate these
 91 quantities, then evaluating forecasts based on the absolute distance between forecast and observed value
 92 forecasted and observed incidence values penalises underprediction (of the reproduction number or growth
 93 rate) less than overprediction by the same amount. This is because for exponential processes errors on the
 94 observed value grow exponentially with the error on the estimated reproduction number or growth rate. For
 95 illustration, consider an incidence forecast issued at time 0 and referring to time t that misses the correct
 96 average growth rate \bar{r}_t by either $-\epsilon$ or $+\epsilon$. Then the ratio of the resulting absolute errors on the scale of
 97 observed incidences y_t is

$$\frac{|y_0 \exp[(\bar{r}_t - \epsilon) \times t] - y_0 \exp(\bar{r}_t t)|}{|y_0 \exp[(\bar{r}_t + \epsilon) \times t] - y_0 \exp(\bar{r}_t t)|} = \exp(-\epsilon t) < 1. \quad (3)$$

99 If one is to measure the ability of forecasts to assess and to forecast the underlying infection dynamics, it
 100 may thus be more desirable to evaluate errors on the scale of the growth rate directly.

101 Evaluating forecasts using the CRPS or WIS means that scores represent a measure of absolute errors.
 102 However, Another argument against using notions of absolute distance between predicted and observed
 103 incidence values is that forecast consumers may find errors on a relative scale easier to interpret and more
 104 useful in order to track predictive performance across targets of different orders of magnitude. Bolin and
 105 Wallin (2023) have proposed the scaled CRPS (SCRPS) which is locally scale invariant; however, it does not
 106 correspond to a relative error measure and lacks a straightforward interpretation as available for the CRPS.

107 A closely related aspect to relative scores (as opposed to absolute scores) is that in the evaluation one may
 108 wish to give similar weight to all considered forecast targets. Lastly, it may be considered desirable to give all
 109 forecast targets similar weight in an overall performance evaluation. As the CRPS typically scales with the
 110 order of magnitude of the quantity to be predicted, this is not the case for the CRPS, which will typically
 111 assign higher scores to forecast targets with high expected values (e.g., in large locations or around the peak
 112 of an epidemic). Bracher et al. (2021a) have argued that this is a desirable feature, directing attention to
 113 situations of particular public health relevance. An evaluation based on absolute errors, however, will assign
 114 little weight to other potentially important aspects, such as the ability to correctly predict future upswings
 115 while observed numbers are still low.

116 In many fields, it is common practice to forecast transformed quantities (see e.g. Taylor (1999) in finance,
 117 Mayr and Ulbricht (2015) in macroeconomics, Löwe et al. (2014) in hydrology or Fuglstad et al. (2015)
 118 in meteorology). While the goal of the transformations is usually often to improve the accuracy of the
 119 predictions, they can also be used to enhance and complement the evaluation process. In this paper, we
 120 argue that the aforementioned issues with evaluating epidemic forecasts based on measures of absolute
 121 error on the natural scale can be addressed by transforming the forecasts and observations prior to scoring
 122 using some strictly monotonic transformation. Strictly monotonic transformations can shift the focus of
 123 the evaluation in a way that may be more appropriate for epidemiological forecasts, while preserving the
 124 propriety of the score guaranteeing that the score remains proper. Many different transformations may be
 125 appropriate and useful, depending on the exact context, the desired focus of the evaluation, and specific
 126 aspects of the forecasts that forecast consumers care most about (see a broader discussion in Section 4).

127 For conceptual clarity and to allow for a more in-depth discussion, we focus mostly on the natural log-
 128 arithm as a particular transformation (referred to as the log transformation in the remainder of this
 129 manuscript) in the context of epidemic phenomena. The particularly attractive transformation in the context
 130 of epidemic phenomena. We refer to this transformation as 'log-transformation' and to scores that have
 131 been computed from log-transformed forecasts and observations as scores 'on the log scale' (as opposed to
 132 scores 'on the natural scale', which involve no transformation). In the theoretical discussion in Section 2,
 133 'log-transformation' and 'log scale' generally refer to a transformation of $\log_e(x)$. For practical applications
 134 (Section 3) we also use these terms to describe a transformation of $\log_e(x + a)$ with a small $a > 0$ in order to
 135 keep the terminology and notation simple. For a prediction target with strictly positive support, the CRPS
 136 after applying a log-transformation can be computed as follows: is given by

$$\begin{aligned}
 \text{CRPS}(F_{\log}, \log y) &= \int_{-\infty}^{\infty} (F_{\log}(x) - 1(x \geq \log y))^2 dx, \\
 \text{CRPS}(F_{\log}, \log y) &= \int_{-\infty}^{\infty} (F_{\log}(x) - 1(x \geq \log y))^2 dx. \tag{4}
 \end{aligned}$$

140 where Here, y is again the observed outcome and F_{\log} is the log-transformed predictive distribution predictive
 141 CDF of the log-transformed outcome, i.e.,

$$F_{\log}(x) = F(\exp(x)), \tag{5}$$

143 with F the CDF on the original scale. Instead of a score representing the magnitude of absolute errors,
 144 applying a log-transformation prior to the CRPS yields a score which a) measures relative error (see Section
 145 2.1), b) provides a measure for how well a forecast captures the exponential growth rate of the target quantity
 146 (see Section 2.2) and c) is less dependent on the expected order of magnitude of the quantity to be predicted
 147 (see Section 2.3). We therefore argue that such evaluations on the logarithmic scale should complement the
 148 prevailing evaluations on the natural scale. Other transformations may likewise be of interest. We briefly
 149 explore the square root transformation as an alternative transformation. Our analysis mostly focuses on
 150 the CRPS (or WIS) as an evaluation metric for probabilistic forecasts, given its widespread use throughout
 151 the COVID-19 pandemic. We note that the logarithmic score has scale invariance properties which imply
 152 that score differences between different forecasts are invariant to strictly monotonic transformations (see
 153 Lehmann 1950 on corresponding properties of likelihood ratios and Diks et al. 2011). The question of the
 154 right scale to evaluate forecasts on does therefore not arise for the log score.

155 The remainder of the article is structured as follows. In Sections 2.1–2.3 we provide some mathematical
 156 intuition on applying the log-transformation prior to evaluating the CRPS, highlighting the connections to
 157 relative error measures, the epidemic growth rate and variance stabilizing transformations. We then discuss
 158 the effect of the log-transformation on forecast rankings (Section 2.4) as well as practical considerations for
 159 applying transformations in general and the log-transformation in particular (Section 2.5) and the effect of
 160 the log-transformation on forecast rankings (Section 2.4). To analyse the real-world implications of the log-
 161 transformation we use forecasts submitted to the European COVID-19 Forecast Hub (European Covid-19

¹⁶² Forecast Hub, 2021; Sherratt et al., 2022, Section 3). Finally, we provide scoring recommendations, discuss
¹⁶³ alternative transformations that may be useful in different contexts, and suggest further research avenues
¹⁶⁴ (Section 4).

¹⁶⁵ 2 Logarithmic transformation of forecasts and observations

¹⁶⁶ 2.1 Interpretation as a relative error

¹⁶⁷ To illustrate the effect of applying the natural logarithm prior to evaluating forecasts we consider the absolute
¹⁶⁸ error, which the CRPS and WIS generalize to probabilistic forecasts. We assume strictly positive support
¹⁶⁹ (meaning that no specific handling of zero values is needed), a restriction we will address when applying this
¹⁷⁰ transformation in practice. When considering a point forecast \hat{y} for a quantity of interest y , such that

$$\begin{aligned} \text{171} \quad & y = \hat{y} + \varepsilon, \\ \text{172} \quad & \underline{y} = \hat{y} + \varepsilon, \end{aligned} \tag{6}$$

¹⁷⁴ the absolute error is given by $|\varepsilon|$. When taking the logarithm of the forecast and the observation first, thus
¹⁷⁵ considering

$$\begin{aligned} \text{176} \quad & \log y = \log \hat{y} + \varepsilon^*, \\ \text{177} \quad & \underline{\log y} = \log \hat{y} + \varepsilon^*, \end{aligned} \tag{7}$$

¹⁷⁹ the resulting absolute error $|\varepsilon^*|$ can be interpreted as an approximation of various common relative error
¹⁸⁰ measures. Using that $\log(a) \approx a - 1$ if a is close to 1, we get

$$\begin{aligned} \text{181} \quad & |\varepsilon^*| = |\log \hat{y} - \log y| = \left| \log \left(\frac{\hat{y}}{y} \right) \right| \quad \text{if } \hat{y} \approx y \quad \left| \frac{\hat{y}}{y} - 1 \right| = \left| \frac{\hat{y} - y}{y} \right|. \\ \text{182} \quad & \underline{|\varepsilon^*|} = |\log \hat{y} - \log y| = \left| \log \left(\frac{\hat{y}}{y} \right) \right| \quad \text{if } \hat{y} \approx y \quad \left| \frac{\hat{y}}{y} - 1 \right| = \left| \frac{\hat{y} - y}{y} \right|. \end{aligned} \tag{8}$$

¹⁸³ The absolute error after log transforming is thus an approximation of the *absolute percentage error* (APE,
¹⁸⁴ Gneiting, 2011) as long as forecast and observation are close. As we assumed that $\hat{y} \approx y$, we can also
¹⁸⁵ interpret it as an approximation of the *relative error* (RE)

$$\begin{aligned} \text{186} \quad & \text{(RE, Gneiting, 2011)} \\ \text{187} \quad & \left| \frac{\hat{y} - y}{\hat{y}} \right| \end{aligned} \tag{9}$$

¹⁸⁸ and the *symmetric absolute percentage error* (SAPE)¹⁸⁹

$$\left| \frac{\hat{y} - y}{y/2 + \hat{y}/2} \right|.$$

¹⁸⁹ ; see e.g., Flores 1986)

$$\left| \frac{\hat{y} - y}{y/2 + \hat{y}/2} \right|. \tag{10}$$

191 As Figure 1 shows, the alignment with the SAPE is in fact the closest and holds quite well even if predicted
 192 and observed value differ by a factor of two or three. Generalising to probabilistic forecasts, the CRPS
 193 applied to log-transformed forecasts and outcomes can thus be seen as a probabilistic counterpart to the
 194 symmetric absolute percentage error, which offers an appealing intuitive interpretation.

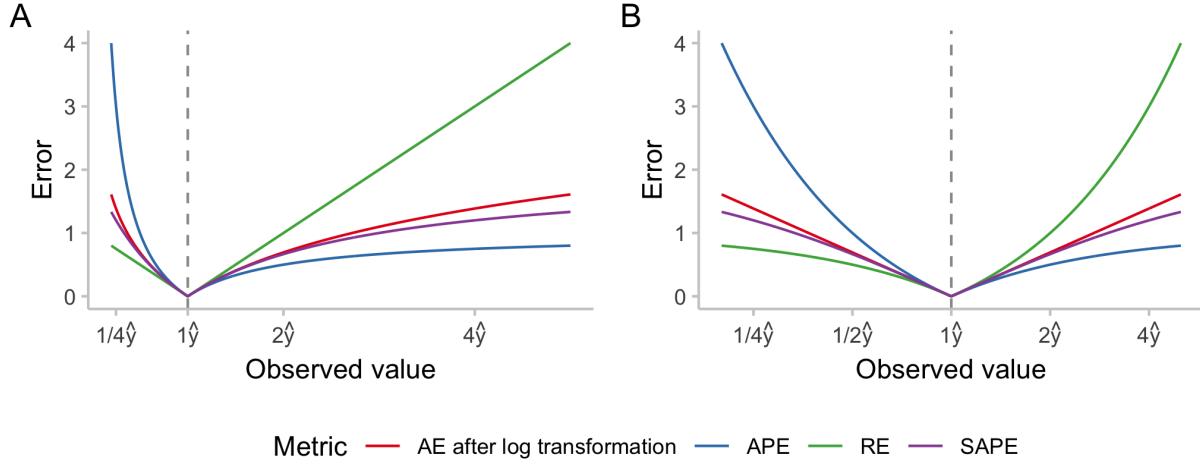


Figure 1: Numerical comparison of different measures of relative error: absolute percentage error (APE), relative error (RE), symmetric absolute percentage error (SAPE) and the absolute error applied to log-transformed predictions and observations. We denote the predicted value by \hat{y} and display errors as a function of the ratio of observed and predicted value. A: x-axis shown on a linear scale. B: x-axis shown on a logarithmic scale.

195 2.2 Interpretation as scoring the exponential growth rate

196 Another interpretation for the log-transform is possible if the generative process is framed as exponential
 197 with a time-varying growth rate $r(t)$ (see, e.g., Wallinga and Lipsitch, 2007), i.e.

$$198 \quad \frac{d}{dt}y(t) = r(t)y(t) \\ 199 \quad \underline{\underline{\frac{d}{dt}y(t)}} = \underline{\underline{r(t)y(t)}} \quad (11)$$

201 which is solved by

$$202 \quad \underline{\underline{y(t) = y_0 \exp\left(\int_0^t r(t')dt'\right)}} = y_0 \exp(\bar{r}t) \\ 203 \quad \underline{\underline{y(t) = y_0 \exp\left(\int_0^t r(t')dt'\right)}} = y_0 \exp(\bar{r}_t t) \quad (12)$$

205 where y_0 is an initial data point and \bar{r}_t is the mean of the growth rate between the initial time point 0 and
 206 time t .

207 If a forecast $\hat{y}(t)$ for the value of the time series at time t is issued at time 0 based on the data point y_0 then
 208 the absolute error after log transformation is

$$209 \quad \underline{\epsilon^* = |\log[\hat{y}(t)] - \log[y(t)]|}$$

210 $= |\log [y_0 \exp(\bar{r}t)] - \log [y_0 \exp(\bar{r}t)]|$
 211 $= t |\bar{r} - \bar{r}|$
 212
 213 $\epsilon^* = |\log [\hat{y}(t)] - \log [y(t)]|$
 $= |\log [y_0 \exp(\hat{r}_t t)] - \log [y_0 \exp(\bar{r}_t t)]|$
 $= t |\hat{r}_t - \bar{r}_t|$

(13)

214 where \hat{r}_t is the true mean growth rate and \hat{r}_t is the forecast mean growth rate. We thus evaluate the error
 215 in the mean exponential growth rate, scaled by the length of the time period considered. Again generalising
 216 this to the CRPS and WIS implies a probabilistic evaluation of forecasts of the epidemic growth rate.

217 2.3 Interpretation as a variance-stabilising transformation

218 When evaluating models across sets of forecasting tasks, it may be desirable for each target to have a similar
 219 impact on the overall results. In disease incidence forecasting, this is not the case when using the CRPS
 220 This could be motivated by the assumption that forecasts from different geographical units and time periods
 221 provide similar amounts of information about how well a forecaster performs. One would then like the
 222 resulting scores to be independent of the order of magnitude of the target to predict. CRPS values on
 223 the natural scale, as the latter typically scales however, typically scale with the order of magnitude of the
 224 quantity to be predicted. Average scores are then dominated by the results achieved for targets with high
 225 expected outcomes in a way that does not necessarily reflect the underlying predictive ability well.

226 Specifically, if If the predictive distribution for the quantity Y equals the true data-generating process F (an
 227 ideal forecast), the expected CRPS is given by (Gneiting and Raftery, 2007)

228 $\mathbb{E}[\text{CRPS}(F, y)] = 0.5 \times \mathbb{E}|Y - Y'|,$
 229 $\mathbb{E}[\text{CRPS}(F, y)] = 0.5 \times \mathbb{E}|Y - Y'|,$

(14)

230 where Y and Y' are independent samples from F . This corresponds to half the *mean absolute difference*,
 231 which is a measure of dispersion. If F is well-approximated by a normal distribution $N(\mu, \sigma^2)$, the approxi-
 232 mation

233 $\mathbb{E}_F[\text{CRPS}(F, y)] \approx \frac{\sigma}{\sqrt{\pi}}$
 234 $\mathbb{E}_F[\text{CRPS}(F, y)] \approx \frac{\sigma}{\sqrt{\pi}}$

(15)

235 can be used. This means that the expected CRPS scales roughly with the standard deviation, which
 236 in turn typically increases with the mean in epidemiological forecasting. In order to make the expected
 237 CRPS independent of the expected outcome, a *variance-stabilising transformation* (VST, Bartlett, 1936)
 238 (VST, Bartlett, 1936; Dunn and Smyth, 2018) can be employed. The choice of this transformation depends
 239 on the mean-variance relationship of the underlying process.

240 If the mean-variance relationship of the data-generating distribution is quadratic with $\sigma^2 = c \times \mu^2$, the
 241 natural logarithm can serve as the VST(Guerrero, 1993). Denoting by F_{\log} the predictive distribution for
 242 $\log(Y)$, we can use the delta method (a first-order Taylor approximation, see e.g., Dunn and Smyth 2018),
 243 to show that

244 $\mathbb{E}_F[\text{CRPS}\{F_{\log}, \log(y)\}] \approx \frac{\sigma/\mu}{\sqrt{\pi}} = \frac{\sqrt{c}}{\sqrt{\pi}}.$

244

$$\mathbb{E}_F[\text{CRPS}\{F_{\log}, \log(y)\}] \approx \frac{\sigma/\mu}{\sqrt{\pi}} = \frac{\sqrt{c}}{\sqrt{\pi}}. \quad (16)$$

245 As σ and μ are linked through the quadratic mean-variance relationship (or linear mean-standard deviation
 246 relationship, $\sigma = \sqrt{c} \times \mu$), the expected CRPS thus stays constant regardless of the expected value of the
 247 data-generating distribution μ . The assumption of a quadratic mean-variance relationship is closely linked
 248 to the aspects discussed in Sections 2.1 and 2.2. It implies that relative errors have constant variance and
 249 can thus be meaningfully compared across different targets. Also, it arises naturally if we assume that our
 250 capacity to predict the epidemic growth rate does not depend on the expected outcome, i.e. does not depend
 251 on the current phase of the epidemic or the order of magnitude of current observations.

252 If the variance-mean-variance relationship is linear with $\sigma^2 = c \times \mu$, as with a Poisson-distributed variable,
 253 the square root is known to be a VST (Dunn and Smyth, 2018). Denoting by $F_{\sqrt{\cdot}}$ the predictive distribution
 254 for \sqrt{Y} , the delta method can again be used to show that

$$\mathbb{E}_F[\text{CRPS}\{F_{\sqrt{\cdot}}, \sqrt{y}\}] \approx \frac{\sigma/\sqrt{\mu}}{2\sqrt{\pi}} = \frac{\sqrt{c}}{2\sqrt{\pi}}.$$

255

$$\mathbb{E}_F[\text{CRPS}\{F_{\sqrt{\cdot}}, \sqrt{y}\}] \approx \frac{\sigma/\sqrt{\mu}}{2\sqrt{\pi}} = \frac{\sqrt{c}}{2\sqrt{\pi}}. \quad (17)$$

256 We note that while standard in the derivation of variance-stabilizing transformations, the application of the
 257 delta method in equations (16) and (17) requires the probability mass of F to be tightly distributed. If this
 258 is not the case, the approximation and thus the variance stabilization may be less accurate.

259 To strengthen our intuition on how transforming outcomes prior to applying the CRPS shifts the emphasis
 260 between targets with high and low expected outcomes, Figure 2 shows the expected CRPS of ideal forecasters
 261 under different mean-variance relationships and transformations. We consider a Poisson distribution where
 262 $\sigma^2 = \mu$, a negative binomial distribution with size parameter $\theta = 10$ and thus $\sigma^2 = \mu + \mu^2/10$, and a
 263 truncated normal distribution with practically constant variance. We see that when applying the CRPS on
 264 the natural scale, the expected CRPS grows with monotonically as the variance of the predictive distribution
 265 (which is equal to the data-generating distribution for the ideal forecaster) increases. The expected CRPS is
 266 constant only for the distribution with constant variance, and grows in μ for the other two. When applying a
 267 log-transformation first, the expected CRPS is almost independent of μ for the negative binomial distribution
 268 and large μ , while smaller targets have higher expected CRPS in case of the Poisson distribution and the
 269 normal distribution with constant variance. When applying a square-root-transformation before the CRPS,
 270 the expected CRPS is independent of the mean for the Poisson-distribution, but not for the other two (with
 271 a positive relationship in the normal case and a negative one for the negative binomial). As can be seen in
 272 Figures 2 and SI.3, the approximations presented above in equations (16) and (17) work quite well for our
 273 simulated example.

274 **2.4 Practical considerations** Effects on model rankings

275 Transformations that are strictly monotonic are permissible in Rankings between different forecasters based
 276 on the CRPS may change when making use of a transformation, both in terms of aggregate and individual
 277 scores. We illustrate this in Figure 3 with two forecasters, A and B, issuing two different distributions
 278 with different dispersion. When showing the obtained CRPS as a function of the observed value, it can be
 279 seen that the sense that they maintain the propriety of the score. This is because even though rankings
 280 of models may change forecasts will in expectation still minimise their score if they report a predictive
 281 distribution that is equal to the data-generating distribution. This condition holds for both the log and
 282 square root transformations, as well as many others. However, the order of the operations matters, and

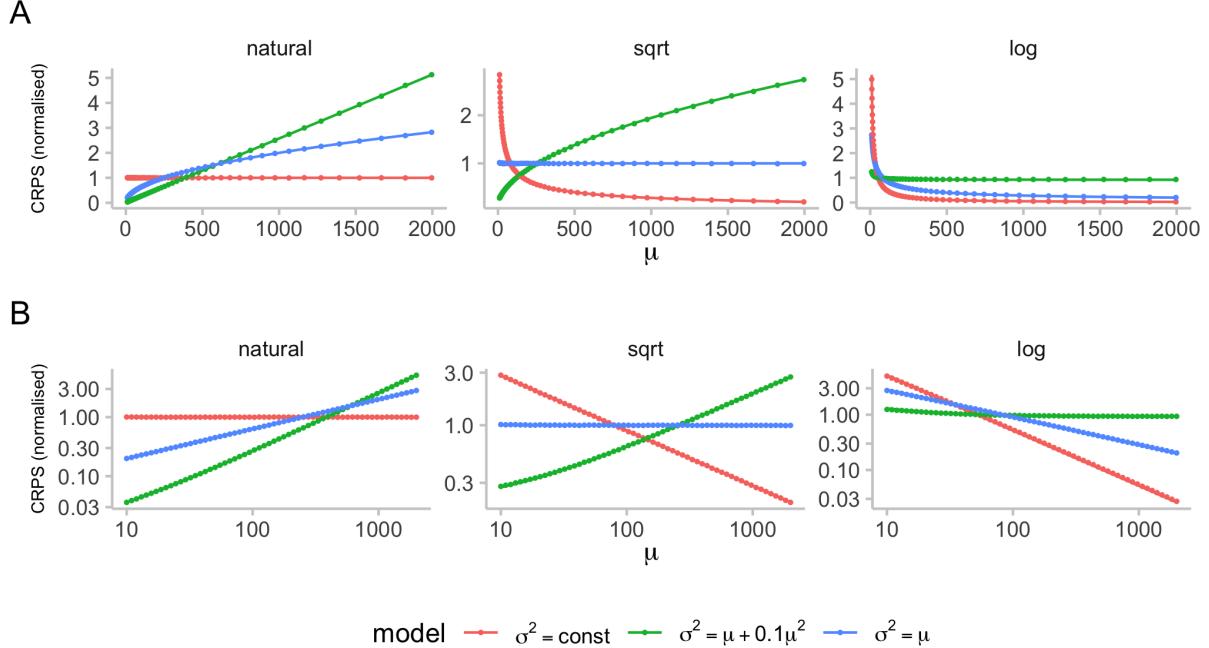


Figure 2: Expected CRPS scores as a function of the mean and variance of the forecast quantity. We computed expected CRPS values for three different distributions, assuming an ideal forecaster with predictive distribution equal to the [true underlying](#) (data-generating) distribution. These expected CRPS values [where were](#) computed for different predictive means based on 10,000 samples each and are represented by dots. Solid lines show the corresponding [approximation approximations](#) of the expected CRPS [based on an assumed normal distribution as discussed in section 2.3 from equations \(16\) and \(17\)](#). Figure SI.3 shows the quality of the approximation in more detail. The first distribution (red) is a truncated normal distribution with constant variance (we chose $\sigma = 1$ in order to only obtain positive samples). The second (green) is a negative binomial distribution with variance $\theta = 10$ and variance $\sigma^2 = \mu + 0.1\mu^2$. The third (blue) is a Poisson distribution with $\sigma^2 = \mu$. To make the scores for the different distributions comparable, scores were normalised to one, meaning that the mean score for every distribution (red, green, blue) is one. A: Normalised expected CRPS for ideal forecasts with increasing means for three distribution with different relationships between mean and variance. Expected CRPS was computed on the natural scale (left), after applying a square-root transformation (middle), and after adding one and applying a log-transformation to the data (right). B: A but with x [axis and y axes](#) on the log scale.

284 applying a transformation after scores have been computed generally does not guarantee propriety. In the
 285 ease of log transforms, taking the logarithm of the scores , rather than scoring the log-transformed forecasts
 286 and data, results in an improper score. This is because taking the logarithm of the CRPS (or WIS) results
 287 in a score that does not penalise outliers enough and therefore incentivises overconfident predictions. We
 288 illustrate this point using simulated data in Figure SI.1, where it can easily be seen that overconfident
 289 models perform best ranking between the two forecasters may change when scoring the forecast on the
 290 logarithmic, rather than the natural scale. In particular, on the natural scale, forecaster A, who issues a
 291 more uncertain distribution, receives a better score than forecaster B for observed values far away from the
 292 centre of the respective predictive distribution. On the log scale, however, forecaster A receives a lower score
 293 for large observed values, being more heavily penalised for assigning large probability to small values (which,
 294 in relative terms, are far away from the actual observation). We note that the chosen example involving a
 295 geometric forecast distribution is somewhat constructed; as shown in Section 3.4 and Figure 8A, rankings
 296 between models in practice stay quite stable for a single forecast.

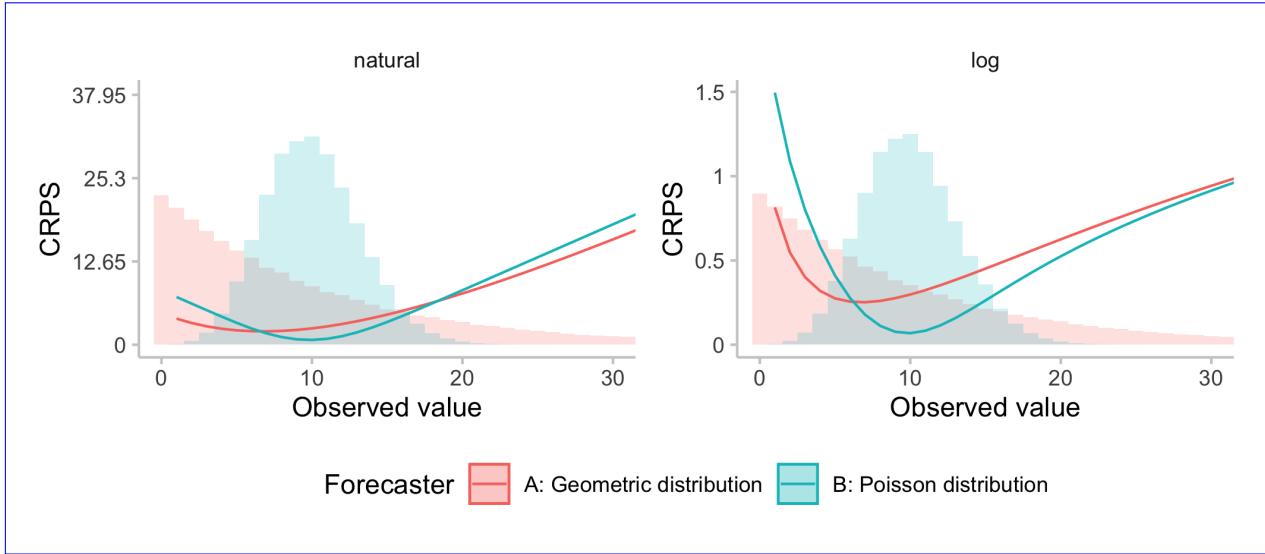


Figure 3: Illustration of the effect of the log-transformation of the ranking for a single forecast. Shown are CRPS (or WIS, respectively) values as a function of the observed value for two forecasters. Model A issues a geometric distribution (a negative binomial distribution with size parameter $\theta = 1$) with mean $\mu = 10$ and variance $\sigma^2 = \mu + \mu^2 = 110$, while Model B issues a Poisson distribution with mean and variance equal to 10. Zeroes in this illustrative example were handled by adding one before applying the natural logarithm.

297 Overall model rankings would be expected to differ more when scores are averaged across multiple forecasts
 298 or targets. The change in rankings of aggregate scores usually is mainly driven by the order of magnitude of
 299 scores for different forecast targets across time, location and target type and less so by the kind of changes in
 300 model rankings for single forecasts discussed above (see Figure 8 for a practical example). Large observations
 301 will dominate average CRPS values when evaluation is done on the natural scale, but much less so after log
 302 transformation. Depending on how different models perform across targets of different orders of magnitude,
 303 rankings in terms of the log-WIS average scores may change when applying a transformation.

304 2.5 Practical considerations and other transformations

305 In practice, one issue with the log transform is that they are it is not readily applicable to negative numbers
 306 or zero values, which need to be removed or otherwise handled. One common approach to deal with zeros
 307 this end is to add a small positive quantity, such as $a = 1$, to all observations and predictions before taking
 308 the logarithm (Bellégo et al., 2022). This still represents a strictly monotonic transformation and therefore
 309 preserves the propriety of the resulting score. The choice of the quantity to add does however influence
 310 but the choice of a does influence scores and rankings, as (measures of relative errors shrink when adding a
 311 constant the larger the chosen value a to the forecast and the observation. We illustrate this in Figure SI.2.
 312). As a rule of thumb, if if $x > 5a$, the difference between $\log(x + a)$ and $\log(x)$ is small, and it becomes
 313 negligible if $x > 50a$. Choosing a suitable offset a thus balances two competing concerns: on the one hand,
 314 choosing a small a makes sure that the transformation is as close to a natural logarithm as possible and
 315 scores can be interpreted as outlined above in the previous sections. On the other hand, choosing a larger
 316 a can help stabilise scores for forecasts and observations close to zero, avoiding giving excessive weight to
 317 forecasts for small quantities of small quantities. For increasing a , less relative weight is given to smaller
 318 forecast targets. For very large values of a , $\log(x + a)$ is roughly linear in x , so that using a very large a
 319 implies similar relative weighting as applying no transformation at all. In practice, a user could explore the
 320 effect of different values of a graphically and choose a such that the relative weightings of times and regions
 321 with high and low incidence correspond to their preferences (see Figure 6 in our example application, Section

322 3).

323 A related issue occurs when the predictive distribution has a large probability mass on zero (or on very small
324 values), as this can translate into an excessively wide forecast in relative terms. ~~This can be seen In our~~
325 ~~applied example this is illustrated~~ in Figure SI.7. ~~HereIn such instances~~, the dispersion component of the
326 WIS is inflated for scores obtained after applying the natural logarithm because forecasts contained zero in
327 its prediction intervals.

328 2.6 Effects on model rankings

329 Rankings between different forecasters based on the CRPS may change when making use of a transformation,
330 both in terms of aggregate and individual scores . We illustrate this in Figure 3 with two forecasters, A
331 and B, issuing two different distributions with different dispersion. When showing the obtained CRPS as a
332 function of the observed value, it can be seen that the ranking between the two forecasters may change when
333 scoring the forecast on the logarithmic, rather than the natural scale. In particular, on the natural scale,
334 forecaster A, who issues a more uncertain distribution, receives a better score than forecaster B for observed
335 values far away from the centre of To deal with this issue one could choose to use a higher a value when
336 applying a transformation $\log(x + a)$, for example $a = 10$ instead of the respective predictive distribution.
337 On the log scale, however, forecaster A receives a lower score for large observed values, being more heavily
338 penalised for assigning large probability to small values (which, in relative terms, are far away from the
339 actual observation) . $a = 1$ that we chose to use.

340 Illustration of the effect of the log-transformation of the ranking for a single forecast. Shown are CRPS
341 (or WIS, respectively) values as a function of the observed value for two forecasters. Model A issues a
342 geometric distribution (a negative binomial distribution with size parameter $\theta = 1$) with mean $\mu = 10$ and
343 variance $\sigma^2 = \mu + \mu^2 = 110$, while Model B issues a Poisson distribution with mean and variance equal to
344 10. Zeroes in this illustrative example were handled by adding one before applying the natural logarithm.
345 A natural question is which other transformations could be applied and whether resulting scores remain
346 (strictly) proper. In principle, any transformation function can be applied simultaneously to forecasts and
347 observations as long as the definition of the transformation is independent of the forecasts and any quantities
348 unknown at the time of forecasting, including the observed value. This simply corresponds to a re-definition
349 of the forecasting target. However, applying non-invertible transformations leads to a loss in information
350 conveyed by forecasts, which we consider undesirable. The resulting score will be proper, but it may not be
351 strictly proper anymore (as forecasts differing from the forecaster's true belief on the original scale may be
352 identical on the transformed scale). When using the CRPS or the WIS, it seems most appropriate to use
353 only strictly monotonic transformations such as the natural logarithm or the square root as otherwise the
354 encoded notion of distance may become meaningless.

355 Overall model rankings would be expected to differ even more when scores are averaged across multiple
356 forecasts or targets. The change in rankings of aggregate scores is mainly driven by Some other strictly
357 monotonic transformations that can be applied are scaling by the population size or scaling by past observations.
358 The latter, as discussed in Section 4, is similar to applying a log-transformation, but corresponds to evaluating
359 a forecast of multiplicative, rather than exponential growth rates. The arising issue of dividing by zero can
360 again be solved by adding a small offset a . Scaling a forecast by the later observed value (as opposed to scaling
361 by past observations) is generally not permissible as it can result in improper scores (see Lerch et al. 2015
362 on the closely related topic of weighting scores with a function of the observed value). Similarly, scaling
363 forecasts and observations by a function of the predictive distribution (like the predictive mean) may lead
364 to improper scores; however, we are unaware of existing theoretical arguments on this.

365 When applying a transformation, the order of magnitude of scores for different forecast targets across time,
366 location and target type and less so by the kind of changes in model rankings for single forecasts discussed
367 above. Large observations will dominate average CRPS values when evaluation is done on the natural scale,

368 but much less so after log transformation. Depending on the relationship between the mean and variance
369 of the forecast target, a log-transformation may even lead to systematically larger scores assigned to small
370 forecast targets, as illustrated in Figure 2. the operations matters, and applying a transformation after
371 scores have been computed generally does not guarantee that the score remains proper. In the case of log
372 transforms, taking the logarithm of the CRPS values, rather than scoring the log-transformed forecasts and
373 data, results in an improper score. We illustrate this point using simulated data in Figure SI.1, where it
374 can be seen that in the example overconfident models perform best in terms of the log WIS. We note that
375 strictly speaking, re-scaling average scores by the average score of a baseline model or average scores across
376 different models to obtain skill scores likewise leads to improper scores (Gneiting and Raftery, 2007). The
377 application of such skill scores, however, is established practice and considered largely unproblematic.

378 We note that in the practical evaluation of operational forecasting systems several additional challenges
379 arise, which we do not study in detail. These concern e.g., the removal of outlying observations and forecasts
380 and the handling of missing forecasts. The solutions we employed in practice are provided in Section 3.1.

381 3 Empirical example: the European Forecast Hub

382 3.1 Setting

383 As an empirical comparison of evaluating forecasts on the natural and on the log scale, we use forecasts from
384 the European Forecast Hub (European Covid-19 Forecast Hub, 2021; Sherratt et al., 2022). The European
385 COVID-19 Forecast Hub is one of several COVID-19 Forecast Hubs (Cramer et al., 2021; Bracher et al.,
386 2021b) which have been systematically collecting, aggregating and evaluating forecasts of several COVID-19
387 targets created by different teams every week. Forecasts are made one to four weeks ahead into the future
388 and follow a quantile-based format with a set of 23 quantiles (0.01, 0.025, 0.05, ..., 0.5, ...0.95, 0.975, 0.99).

389 The forecasts used for the purpose of this illustration are forecasts submitted between the 8th of March
390 2021 and the 5th of December 2022 for reported cases and deaths from COVID-19. Target dates range from
391 the 13th of March 2021 to the 10th of December 2022, for a total of 92 weeks. See Sherratt et al. (2022)
392 for a more thorough description of the data. We filtered all forecasts submitted to the Hub to only include
393 the seven models which have submitted forecasts for both deaths and cases for 4 horizons in 32 locations
394 on at least 46 forecast dates (see Figure SI.4). We removed all observations marked as data anomalies by
395 the European Forecast Hub (Sherratt et al., 2022) as well as all remaining negative observed values. These
396 anomalies made up a relevant fraction of all observations. On average across locations, 12.1 out of 92 (13.2%)
397 observations were removed for cases and 12.4 out of 92 (13.5%) for deaths. Figure SI.5 displays the number
398 of anomalies removed for each location. In addition, we filtered out erroneous forecasts a small number of
399 erroneous forecasts that were in extremely poor agreement with the observed data, as defined by any of the
400 conditions listed in Table SI.2. Those forecasts were removed. Figure SI.6 shows the percentage of forecasts
401 removed for each model. Those few (less than 0.2% of forecasts for each model) erroneous outlier forecasts
402 had excessive influence on average scores and relative skill scores in a way that was not representative of
403 normal model behaviour. We removed them here in order to be better able to better illustrate the effects of
404 the log-transformation on scores and eliminating distortions caused by outlier forecasters. that one would
405 expect in a well-behaved scenario. In a regular forecast evaluation such erroneous forecasts should usually
406 not be removed and would count towards overall model scores.

407 All predictive quantiles were truncated at 0. We applied the log-transformation after adding a constant
408 $a = 1$ to all predictions and observed values. The choice of $a = 1$ in part reflects convention, but also
409 represents a suitable choice as it avoids giving excessive weight to forecasts close to zero, while at the same
410 time ensuring that scores for observations > 5 can be interpreted reasonably. The analysis was conducted
411 in R (R Core Team, 2022), using the `scoringutils` package (Bosse et al., 2022) for forecast evaluation. All

412 code is available on GitHub (<https://github.com/epiforecasts/transformation-forecast-evaluation>). Where
 413 not otherwise stated, we report results for a two-week-ahead forecast horizon.

414 In addition to the WIS we use pairwise comparisons (Cramer et al., 2021) to evaluate the relative performance
 415 of models across countries in the presence of missing forecasts. In the first step, score ratios are computed
 416 for all pairs of models by taking the set of overlapping forecasts between the two models and dividing the
 417 score of one model by the score achieved by the other model. The relative skill for a given model compared
 418 to others is then obtained by taking the geometric mean of all score ratios which involve that model. Low
 419 values are better, and the “average”–“average” model receives a relative skill score of 1.

420 **3.2 Illustration and qualitative observations**

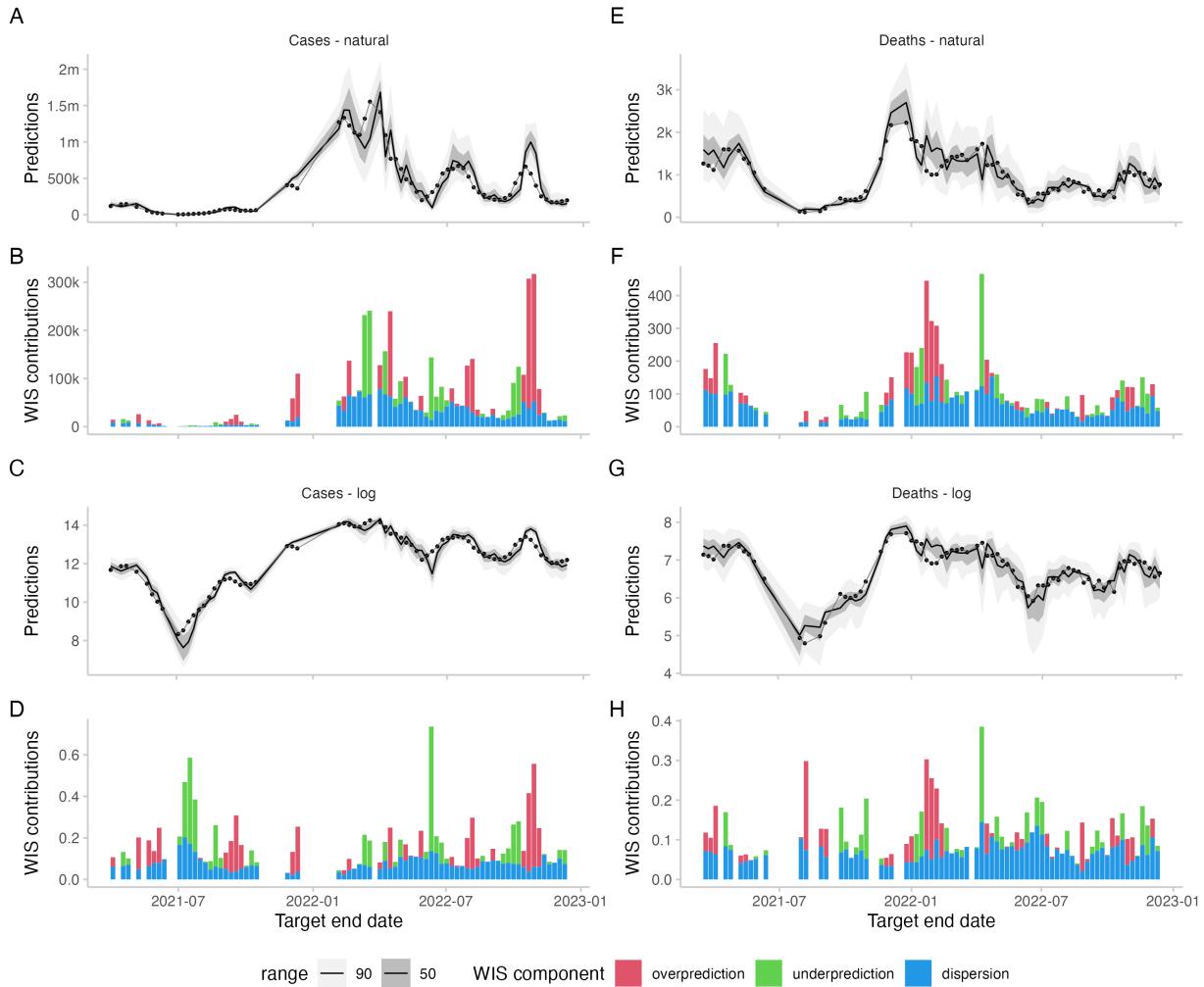


Figure 4: Forecasts and scores for two-week-ahead predictions from the EuroCOVIDhub-ensemble made in Germany. Missing values are due to data anomalies that were removed (see section 3.1. A, E: 50% and 90% prediction intervals and observed values for cases and deaths on the natural scale. B, F: Corresponding scores. C, G: Forecasts and observations on the log scale. D, H: Corresponding scores.

421 When comparing examples of forecasts on the natural scale with those on the log scale (see Figures 4, SI.7,

422 SI.8) a few interesting patterns emerge. Missing the peak, i.e. predicting increasing numbers while actual
 423 observations are already falling, tends to contribute a lot to overall scores on the natural scale (see forecasts
 424 during the peak in May 2022 in May in Figure 4A, B). On the log scale, these have less of an influence,
 425 as errors are smaller in relative terms (see 4C, D). Conversely, failure to predict an upswing while numbers
 426 are still low, is less severely punished-penalised on the natural scale (see forecasts in July 2021 and to a
 427 lesser extent in July 2022 in Figure 4 A, B), as overall absolute errors are low. On the log scale, missing
 428 lower inflection points tends to lead to more severe penalties (see Figure 4C, D)). One can also observe
 429 that on the natural scale, scores tend to track the overall level of the target quantity (compare for example
 430 forecasts for March-July with forecasts for September-October in Figure 4E, F). On the log scale, scores do
 431 not exhibit this behaviour and rather increase whenever forecasts are far away from the truth in relative
 432 terms, regardless of the overall level of observations.

433 Across the dataset, the average number of observed cases and deaths varied considerably by location and
 434 target type (see Figure 5A and B). On the natural scale, scores show a pattern quite similar to the ob-
 435 servations across targets (see Figure5D) and locations (see Figure5C). On the log scale, scores were more
 436 evenly distributed between targets (see Figure5D) and locations (see Figure5C). Both on the natural scale
 437 as well on the log scale, scores increased considerably with increasing forecast horizon (see Figure 5E). This
 438 reflects the increasing difficulty of forecasts further into the future and, for the log scale, corresponds with
 439 our expectations from Section 2.2.

440 To assess the impact of the choice of offset value a we extend the display from Figure 5C by results obtained
 441 under different specifications. Results are shown in Figure 6, where for completeness we also added the
 442 square root transformation. As discussed in Section 2.5, smaller values of a increase the relative weight
 443 of smaller locations in the overall evaluation. In the most extreme considered case $a = 0.001$, the smallest
 444 locations in fact receive the largest weight both for deaths and cases. For very large values (see the third
 445 row of Figure 6), the relative weights strongly resemble those of the evaluation on the natural scale. We
 446 recommend using displays of this type to get an intuition for the role different locations may play for overall
 447 evaluation results.

448 3.3 Regression analysis to determine the variance-stabilizing transformation

449 As argued in Section 2.3, the mean-variance, or mean-CRPS, relationship determines which transformation
 450 can serve as a VST. We can analyse this relationship empirically by running a regression that explains the
 451 WIS (which approximates the CRPS) as a function of the central estimate of the predictive distribution.
 452 We ran the regression

$$453 \log[WIS(F, y)] = \alpha + \beta \times \log[\text{median}(F)], \\ 454 \log[WIS(F, y)] = \alpha + \beta \times \log[\text{median}(F)], \quad (18)$$

455 where the predictive distribution F and the observation y are on the natural scale. This is equivalent to

$$456 \text{WIS}(F, y) = \exp(\alpha) \times \text{median}(F)^\beta, \\ 457 \text{WIS}(F, y) = \exp(\alpha) \times \text{median}(F)^\beta, \quad (19)$$

458 meaning that we estimate a polynomial relationship between the predictive median and achieved WIS. Note
 459 that we are using predictive medians rather than means as only the former are available in the European
 460 COVID-19 Forecast Hub. As ~~the~~-(under the simplifying assumption of normality; see Section 2.3) the
 461 WIS/CRPS of an ideal forecaster scales with the standard deviation(~~see Section 2.3~~), a value of $\beta = 1$
 462 would imply a quadratic median-variance relationship; the natural logarithm could then serve as a VST.
 463 A value of $\beta = 0.5$ would imply a linear median-variance relationship, suggesting the square root as a
 464 VST. We applied the regression to case and death forecasts, ~~pooled across horizons and~~ stratified for one

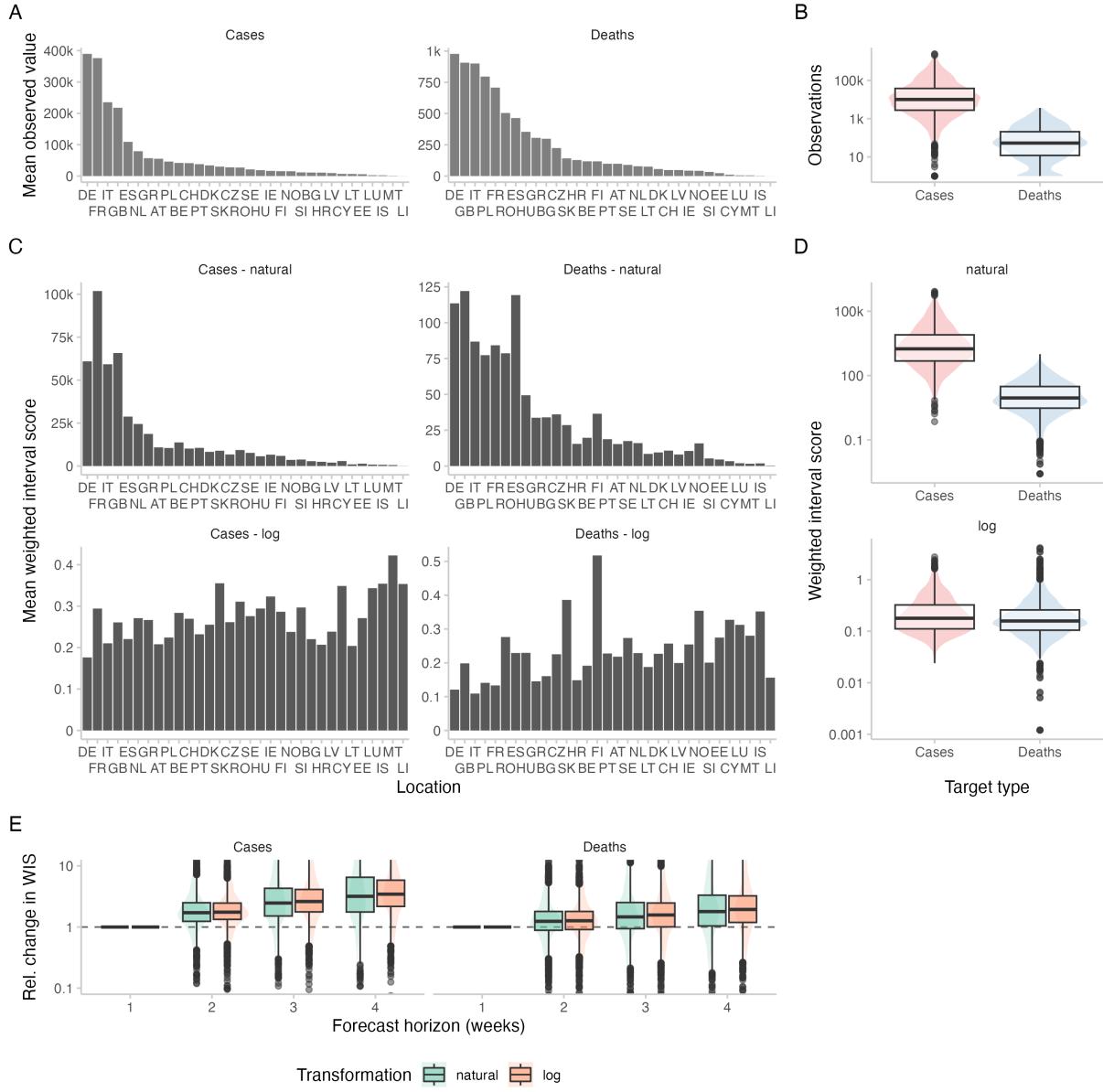


Figure 5: Observations and scores across locations and forecast horizons for the European COVID-19 Forecast Hub data. Locations are sorted according to the mean observed value in that location. A: Average (across all time points) of observed cases and deaths for different locations. B: Corresponding boxplot (y-axis on log-scale) of all cases and deaths. C: Scores for two-week-ahead forecasts from the EuroCOVIDhub-ensemble (averaged across all forecast dates) for different locations, evaluated on the natural scale as well as after transforming counts by adding one and applying the logarithmic scale. D: Corresponding boxplots of all individual scores of the EuroCOVIDhub-ensemble for two-week-ahead predictions. E: Boxplots for the relative change of scores for the EuroCOVIDhub-ensemble across forecast horizons. For any given forecast date and location, forecasts were made for four different forecast horizons, resulting in four scores. All scores were divided by the score for forecast horizon one. To enhance interpretability, the range of visible relative changes in scores (relative to horizon = 1) was restricted to [0.1, 10].

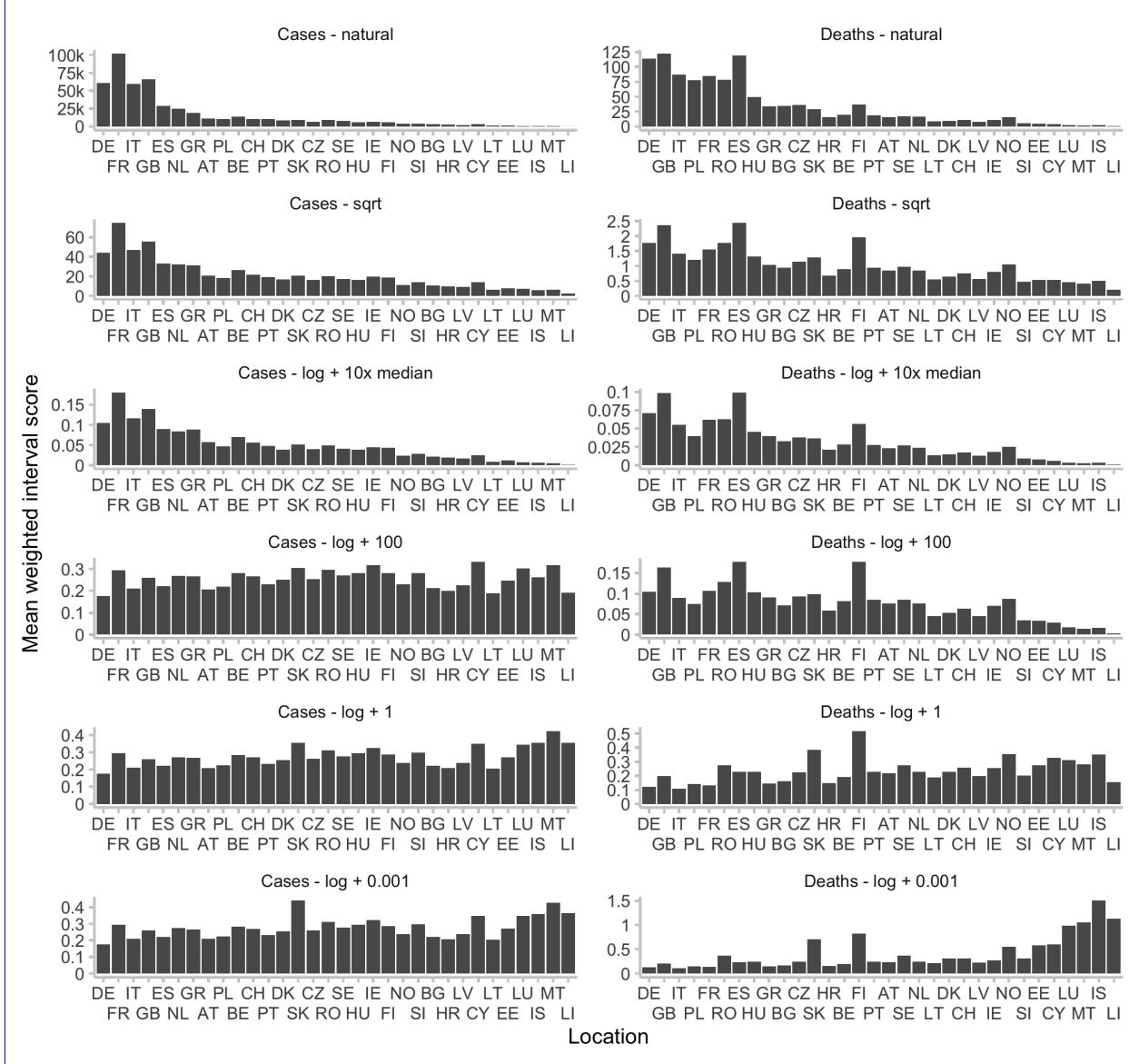


Figure 6: Mean WIS in different locations for different transformations applied before scoring. Locations are sorted according to the mean observed value in that location. Shown are scores for two-week-ahead forecasts of the EuroCOVIDhub-ensemble. On the natural scale (with no transformation prior to applying the WIS), scores correlate strongly with the average number of observed values in a given location. The same is true for scores obtained after applying a square-root transformation, or after applying a log-transformation with a large offset a . For illustrative purposes, a was chosen to be 101630 for cases and 530 for deaths, 10 times the respective median observed value. For large values of a , $\log(x + a)$ grows roughly linearly in x , meaning that we expect to observe the same patterns as in the case with no transformation. For decreasing values of a , we give more relative weight to scores in small locations.

465 through four-week-ahead forecasts. Results are provided in Table 1. It can be seen that the estimates of β
 466 always take a value somewhat below 1, implying a slightly sub-quadratic mean-variance relationship. The
 467 logarithmic transformation should thus approximately stabilize the variance (and WIS), possibly leading to
 468 somewhat higher scores for smaller forecast targets. The square-root transformation, on the other hand, can
 469 be expected to still lead to higher WIS values for targets of higher orders of magnitude.

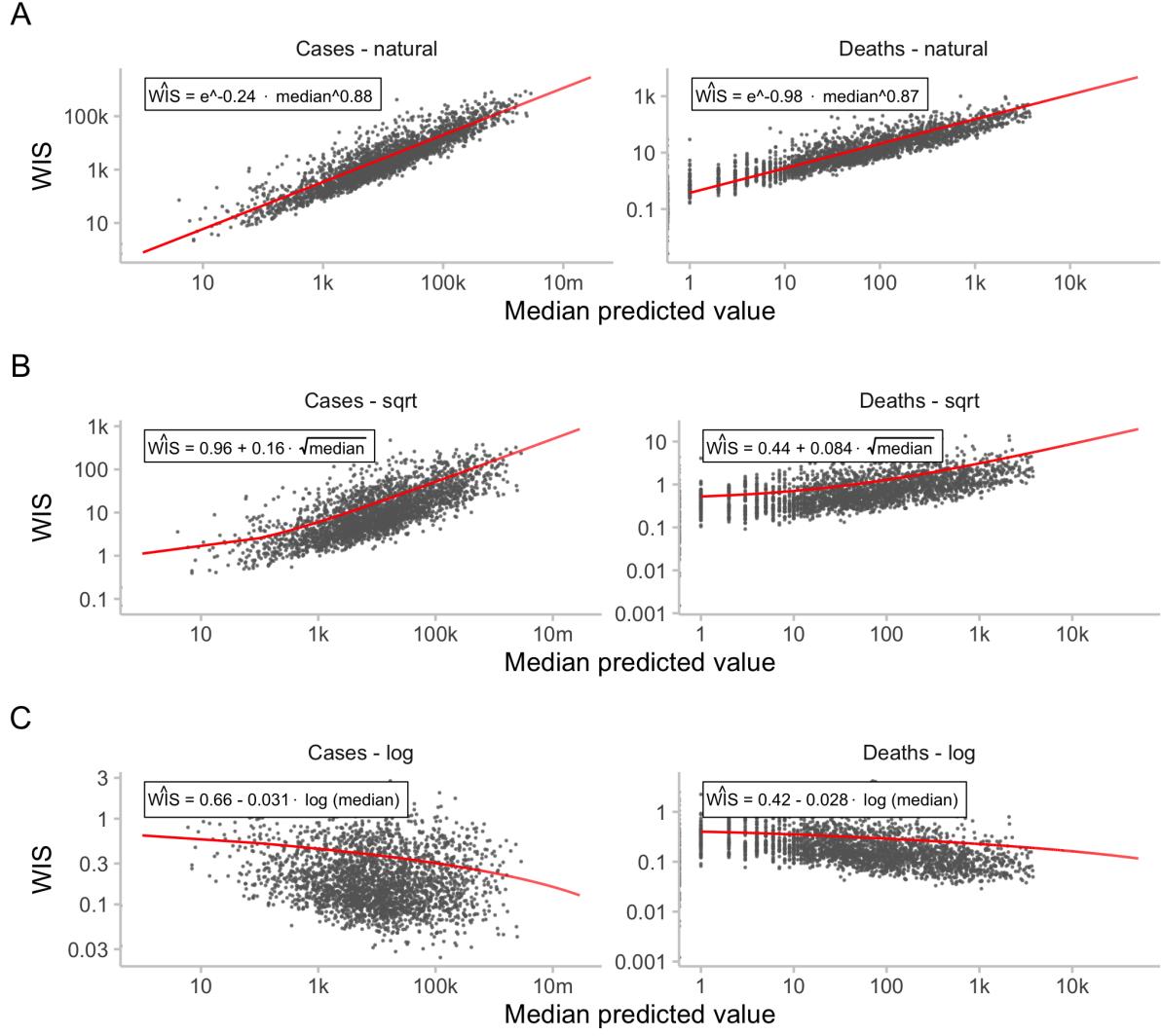


Figure 7: Relationship between median forecasts and scores. Black dots represent WIS values for two-week ahead predictions of the EuroCOVIDhub-ensemble. Shown in red are the regression lines discussed in Section 3.3 shown in Table 1. A: WIS for two-week-ahead predictions of the EuroCOVIDhub-ensemble against median predicted values. B: Same as A, with scores obtained after applying a square-root-transformation to the data. C: Same as A, with scores obtained after applying a log-transformation to the data.

470 To check the relationship after the transformation, we ran the regressions

$$\text{WIS}(F_{\log}, \log y) = \alpha_{\log} + \beta_{\log} \cdot \log(\text{median}(F)),$$

471 $\text{WIS}(F_{\log}, \log y) = \alpha_{\log} + \beta_{\log} \cdot \log(\text{median}(F)),$ (20)

472 where F_{\log} is the predictive distribution for $\log(y)$, and

$$\text{WIS}(F_{\sqrt{y}}, \sqrt{y}) = \alpha_{\sqrt{y}} + \beta_{\sqrt{y}} \cdot \sqrt{\text{median}(F)},$$

473 $\text{WIS}(F_{\sqrt{y}}, \sqrt{y}) = \alpha_{\sqrt{y}} + \beta_{\sqrt{y}} \cdot \sqrt{\text{median}(F)},$ (21)

	Horizon									
all-all	-1.093	-0.963	-0.352	0.201	0.391	0.001	all-Cases	0.036	0.858	0.043
								0.201	0.751	-0.033
									all-Deaths	-0.884
										0.868
										0.273
										0.121
	2	3	4	1	2	3	4			

Table 1: Coefficients of three regressions for the effect of the magnitude of the median forecast on expected scores. The first regression was $\log[WIS(F, y)] = \alpha + \beta \times \log[\text{median}(F)]$, where F is the predictive distribution and y the observed value. The second one was $WIS(F_{\log}, \log y) = \alpha_{\log} + \beta_{\log} \cdot \log(\text{median}(F))$, where F_{\log} is the predictive distribution for $\log y$. The third one was $WIS(F_{\sqrt{\cdot}}, \sqrt{y}) = \alpha_{\sqrt{\cdot}} + \beta_{\sqrt{\cdot}} \cdot \sqrt{(\text{median}(F))}$, $WIS(F_{\sqrt{\cdot}}, \sqrt{y}) = \alpha_{\sqrt{\cdot}} + \beta_{\sqrt{\cdot}} \cdot \sqrt{(\text{median}(F))}$, where $F_{\sqrt{\cdot}}$ is the predictive distribution for \sqrt{y} .

474 where $F_{\sqrt{\cdot}}$ is the predictive distribution on the square-root scale. A value of $\beta_{\log} = 0$ (or $\beta_{\sqrt{\cdot}} = 0$, respectively),
 475 would imply that scores are linearly independent of the median prediction after the transformation. A
 476 value smaller (larger) than 0 would imply that smaller (larger) targets lead to higher scores. As can be seen
 477 from Table 1, the results indeed indicate that small targets lead to larger average WIS when using the log
 478 transform ($\beta_{\log} < 0$), while the opposite is true for the square-root transform ($\beta_{\sqrt{\cdot}} > 0$). The results of
 479 the three regressions are also displayed in Figure 7. In this empirical example, the log transformation thus
 480 helps (albeit not perfectly), to stabilise WIS values, and it does so more successfully than the square-root
 481 transformation. As can be seen from Figure 7, the expected WIS scores for case targets with medians of
 482 10 and 100,000 differ by more than a factor of ten for the square root transformation, but only a factor of
 483 around 2 for the logarithm.

484 3.4 Impact of logarithmic transformation on model rankings

485 For *individual* forecasts, rankings between models for single forecasts are mostly preserved, with differences
 486 increasing across forecast horizons (see Figure 8A). While rankings between forecasters remain similar for
 487 a single forecast, this is not true anymore when looking at rankings obtained after averaging scores across
 488 multiple forecasts made at different times or in different locations. As discussed earlier, scores on the
 489 natural and on the log scale penalise errors very differently, e.g. when looking at performance during peaks
 490 or troughs. When evaluating performance averaged across different forecasts and forecast targets, relative
 491 skill scores of the models therefore change considerably (Figure 8B). The correlation between relative skill
 492 scores also decreases noticeably with increasing forecast horizon.

493 Figure Figure 9 shows the changes in the ranking between different forecasting models. Encouragingly for the
 494 European Forecast Hub, the Hub ensemble, which is the forecast the organisers suggest forecast consumers
 495 make use of, remains the top model across scoring schemes. For cases, the ILM-EKF model and the Forecast
 496 Hub baseline model exhibit the largest change in relative skill scores. For the ILM-EKF model the relative
 497 proportion of the score that is due to overprediction is reduced when applying a log-transformation before
 498 scoring (see Figure 9E). Instances where the model has overshot are penalised less heavily on the log scale,
 499 leading to an overall better score. For the Forecast Hub baseline model, the fact that it often puts relevant
 500 probability mass on zero (see Figure SI.7), leads to worse scores after applying log-transformation due to
 501 large dispersion penalties. For deaths, the baseline model seems to get similarly penalised for its in relative
 502 terms highly dispersed forecasts. The performance of other models changes as well, but patterns are less
 503 discernible on this aggregate level.

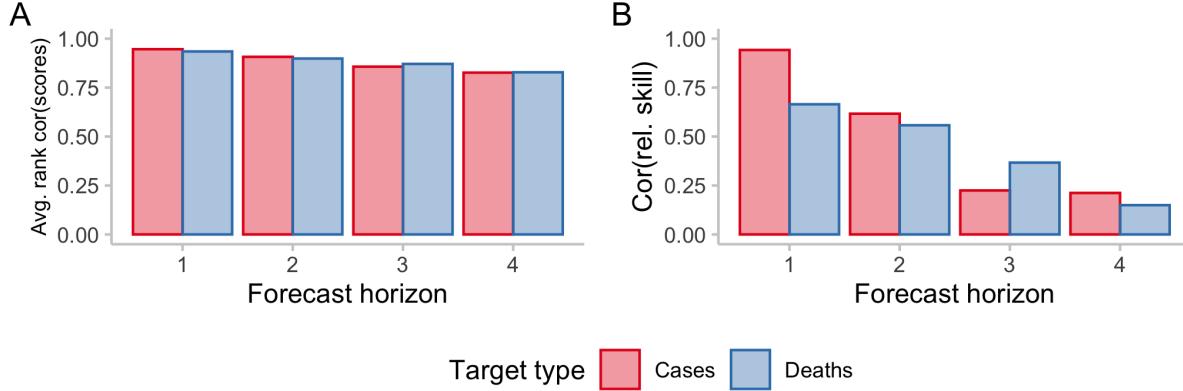


Figure 8: Correlations of rankings on the natural and logarithmic scale. A: Average Spearman rank correlation of scores for individual forecasts from all models. For every individual target (defined by a combination of forecast date, target type, horizon, location), one score was obtained per model. Then, for every forecast target, the Spearman rank correlation was computed between the scores for all models on the natural scale vs. and on the log scale for all the models that had made a forecast for that specific target. All These individual rank correlations were then averaged across locations, and target types time and are displayed stratified by horizon and target types, representing average accordance of model ranks for a single forecast target on the natural and on the log scale. B: Correlation between relative skill scores. For every forecast horizon and target type, a separate relative skill score was computed per model using pairwise comparisons, which is a measure of performance of a model relative to the others for a given horizon and target type that accounts for missing values. The plot shows the correlation between the relative skill scores on the natural vs. on the log scale, representing accordance of overall model performance as judged by scores on the natural and on the log scale.

504 4 Discussion

505 In this paper, we proposed the use of transformations, with a particular focus on the natural logarithmic
 506 transformation, when evaluating forecasts in an epidemiological setting. These transformations can address
 507 issues that arise when evaluating epidemiological forecasts based on measures of absolute error and their
 508 probabilistic generalisations (i.e CRPS and WIS). We showed that scores obtained after log-transforming
 509 both forecasts and observations can be interpreted as a) a measure of relative prediction errors, as well as
 510 b) a score for a forecast of the exponential growth rate of the target quantity and c) as variance stabilising
 511 transformation in some settings. When applying this approach to forecasts from the European COVID-19
 512 Forecast Hub, we found overall scores on the log scale to be more equal across, time, location and target
 513 type (cases, deaths) than scores on the natural scale. Scores on the log scale were much less influenced by
 514 the overall incidence level in a country and showed a slight tendency to be higher in locations with very low
 515 incidences. We found that model rankings changed noticeably.

516 On the natural scale, missing the peak and overshooting was more severely penalised than missing the nadir
 517 and the following upswing in numbers. Both failure modes tended to be more equally penalised on the log
 518 scale (with undershooting receiving slightly higher penalties in our example).

519 Applying a log-transformation prior to the WIS means that forecasts are evaluated in terms of relative
 520 errors and errors on the exponential growth rate, rather than absolute errors. The most important strength
 521 of this approach is that the evaluation better accommodates the exponential nature of the epidemiological
 522 process and the types of errors forecasters who accurately model those processes are expected to make. The
 523 log-transformation also helps avoid issues with scores being strongly influenced by the order of magnitude

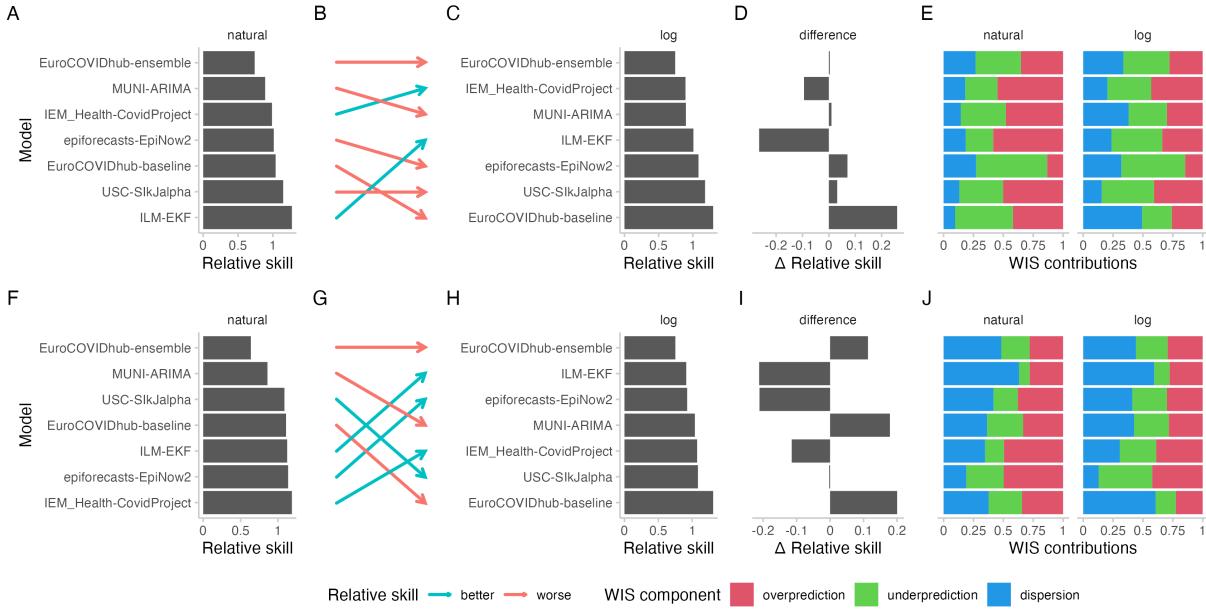


Figure 9: Changes in model ratings as measured by relative skill for two-week-ahead predictions for cases (top row) and deaths (bottom row). A: Relative skill scores for case forecasts from different models submitted to the European COVID-19 Forecast Hub computed on the natural scale. B: Change in rankings as determined by relative skill scores when moving from an evaluation on the natural scale to one on the logarithmic scale. Red arrows indicate that the relative skill score deteriorated when moving from the natural to the log scale, green arrows indicate they improved. C: Relative skill scores based on scores on the log scale. D: Difference in relative skill scores computed on the natural and on the logarithmic scale, ordered as in C. E: Relative contributions of the different WIS components (overprediction, underprediction, and dispersion) to overall model scores on the natural and the logarithmic scale. F, G, H, I, J, K: Analogously for deaths.

of the forecast quantity, which can be an issue when evaluating forecasts on the natural scale. A potential downside is that forecast evaluation is unreliable in situations where observed values are zero or very small. One could argue that this correctly reflect inherent uncertainty about the future course of an epidemic when numbers are small. Users nevertheless need to be aware that this can pose issues in practice. Including very small values in prediction intervals (see e.g. Figure SI.7) can lead to excessive dispersion values on the log scale. Similarly, locations with lower incidences may get disproportionate weight (i.e. high scores) when evaluating forecasts on the log scale. Bracher et al. (2021a) argue that the large weight given it is desirable to give large weight to forecasts for locations with high incidences is a desirable property, as it, as this reflects performance on the targets we should care about most. On the other hand, scoring forecasts on the log scale may be less influenced by outliers and better reflect consistent performance across time, space, and forecast targets. It also gives higher weight to another type of situation one may care about, namely one. Furthermore, decision makers may specifically care about situations in which numbers start to rise from a previously low level.

The log-transformation is only one of many transformations that may be useful and appropriate in an epidemiological context. One obvious option is to apply a population standardization to obtain incidence forecasts e.g., per 100,000 population (Abbott et al., 2022). If one is interested in multiplicative, rather than exponential growth rates, one could convert forecasts into forecasts for the multiplicative growth rate by dividing numbers by the last observed value. We suggested using the natural logarithm as a variance-stabilising transformation (VST) or alternatively the. This is appropriate for variables that are approximately normally distributed and have a quadratic mean-variance relationship with $\sigma^2 = c \times \mu^2$ (this is e.g. approximately true for the negative binomial distribution and large μ). Alternatively, the square-

545 root transformation [can be appropriate](#) in the case of a Poisson distributed variable (Dunn and Smyth, 2018).
546 Other VST like the Box-Cox (Box and Cox, 1964) are conceivable as well. [Another promising transformation would](#)
547 [be to take differences between forecasts on the log scale, or alternatively](#) If one is interested in multiplicative,
548 rather than exponential growth rates, one could instead of applying a log transformation, convert forecasts
549 into forecasts for the multiplicative growth rate by dividing numbers by the last value that was observed at
550 the time the forecast was made. Forecasters would then implicitly predict a separate multiplicative growth
551 rate from today to horizon 1, 2, etc. Instead of dividing by the last observed value, another promising
552 transformation would be to divide each forecast by the forecast of the previous week (and analogously for
553 observations), in order to obtain forecasts for week-to-week growth rates. [One could then also ask forecasters](#)
554 [to provide estimates of the weekly relative change applied to the latest data and subsequent forecast points](#)
555 [directly.](#) [Alternatively, one could also take first differences of values on the log scale.](#) This approach
556 would be akin to evaluating the shape of the predicted trajectory against the shape of the observed trajectory
557 (for a different approach to evaluating the shape of a forecast, see Srivastava et al., 2022). [This Dividing](#)
558 [values by the previous value](#), unfortunately, is not feasible under the current quantile-based format of the
559 Forecast Hubs, as the growth rate of the α -quantile may be different from the α -quantile of the growth-rate.
560 However, it may be an interesting approach if predictive samples are available or if quantiles for weekwise
561 growth rates have been collected. [Potentially, the variance stabilising time-series forecasting literature may](#)
562 [be a useful source of other transformations for various forecast settings.](#)

563 It is possible to go beyond choosing a single transformation by constructing composite scores as a weighted
564 sum of scores based on different transformations. This would make it possible to create custom scores and
565 allow forecast consumers to [choose and](#) assign explicit weights to different qualities of the forecasts they
566 might care about.

567 [In this work, we focused on the CRPS and WIS, which are widely used in the evaluation of epidemic](#)
568 [forecasts. We note that for the logarithmic score, which has also been used e.g., in some editions of the](#)
569 [FluSight challenge Reich et al. \(2019\), the question of the right scale to evaluate forecasts does not arise. It](#)
570 [is known that log score differences between different forecasters are invariant to monotonic transformations](#)
571 [of the outcome variable \(see e.g., Diks et al. 2011\). This is clearly an advantage of the logarithmic score over](#)
572 [the CRPS; however, the logarithmic score is known to have other severe downsides, e.g., its low robustness](#)
573 [to sporadically misguided forecasts; see Bracher et al. \(2021a\) for a more detailed discussion.](#)

574 Exploring transformations is a promising avenue for future work that could help bridge the gap between
575 modellers and policymakers by providing scoring rules that better reflect what forecast consumers care
576 about. [Potentially, the variance stabilising time-series forecasting literature may be a useful source of](#)
577 [transformations for various forecast settings.](#) [In this paper, we did not make any particular assumptions](#)
578 [about policy makers' priorities and preferences.](#) Rather, we aimed to enable users to make an informed
579 choice by showing how different transformations lead to different relative weights for the kinds of prediction
580 errors forecast consumers may care about, such as absolute vs. relative errors or the size of penalties for
581 over- vs. underprediction. [In practice, engagement with decision makers is important to determine what](#)
582 [their priorities are and how different ways to measure predictive importance should be weighed.](#)

583 We have shown that the natural logarithm transformation can lead to significant changes in the relative
584 rankings of models against each other, with potentially important implications for decision-makers who rely
585 on the knowledge of past performance to make a judgement about which forecasts should inform future
586 decisions. While it is commonly accepted that multiple proper scoring rules should usually be considered
587 when comparing forecasts, we think this should be supplemented by considering different transformations of
588 the data to obtain a richer picture of model performance. More work needs to be done to better understand
589 the effects of applying transformations in different contexts, and how they may impact decision-making.

590 **A Supplementary information**

591 **A.1 Alternative Formulation of the WIS**

592 Instead of defining the WIS as an average of scores for individual quantiles, we can define it using an
 593 average of scores for symmetric predictive intervals. For a single prediction interval, the interval **seoren**
 594 **score** (IS) is computed as the sum of three penalty components, dispersion (width of the prediction interval),
 595 underprediction and overprediction,

$$596 \quad \underline{IS_\alpha(F, y)} = (u - l) + \frac{2}{\alpha} \cdot (l - y) \cdot \mathbf{1}(y \leq l) + \frac{2}{\alpha} \cdot (y - u) \cdot \mathbf{1}(y \geq u)$$

597 $\quad \quad \quad = \text{dispersion} + \text{underprediction} + \text{overprediction},$

598

$$599 \quad \underline{IS_\alpha(F, y)} = (u - l) + \frac{2}{\alpha} \cdot (l - y) \cdot \mathbf{1}(y \leq l) + \frac{2}{\alpha} \cdot (y - u) \cdot \mathbf{1}(y \geq u) \quad (22)$$

600 $\quad \quad \quad = \text{dispersion} + \text{underprediction} + \text{overprediction}, \quad (23)$

601 where $\mathbf{1}(\cdot)$ is the indicator function, y is the observed value, and l and u are the $\frac{\alpha}{2}$ and $1 - \frac{\alpha}{2}$ quantiles
 602 of the predictive distribution, i.e. the lower and upper bound of a single central prediction interval. For a
 603 set of K^* prediction intervals and the median m , the WIS is computed as a weighted sum,

$$604 \quad \underline{\text{WIS}} = \frac{1}{K^* + 0.5} \cdot \left(w_0 \cdot |y - m| + \sum_{k=1}^{K^*} w_k \cdot IS_{\alpha_k}(F, y) \right),$$

605

$$606 \quad \underline{\text{WIS}} = \frac{1}{K^* + 0.5} \cdot \left(w_0 \cdot |y - m| + \sum_{k=1}^{K^*} w_k \cdot IS_{\alpha_k}(F, y) \right), \quad (24)$$

607 where w_k is a weight for every interval. Usually, $w_k = \frac{\alpha_k}{2}$ and $w_0 = 0.5$.

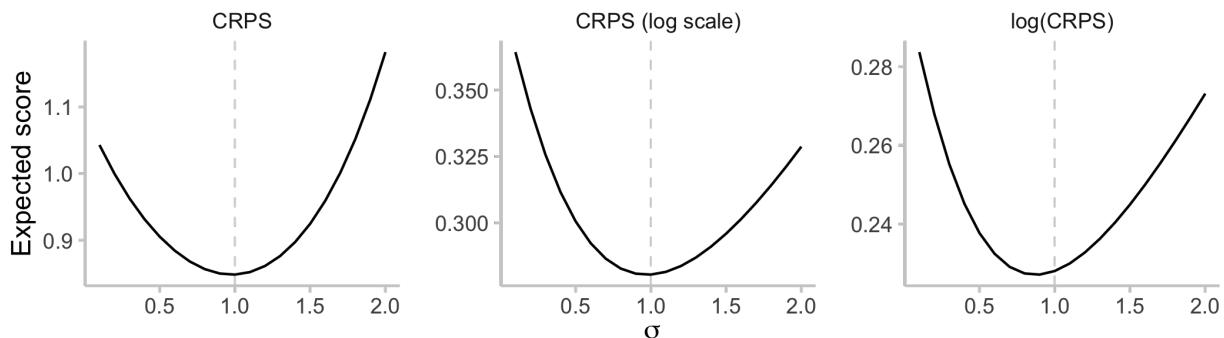


Figure SI.1: Illustration of **impropriety** **the effect** of **log-transformed CRPS** **applying a transformation after scoring**. We assume $Y \sim \text{LogNormal}(0, 1)$ and evaluate the expected CRPS for predictive distributions $\text{LogNormal}(0, \sigma)$ with varying values of $\sigma \in [0.1, 2]$. For the regular CRPS (left) and CRPS applied to log-transformed outcomes (middle), the lowest expectation is achieved for the true value $\sigma = 1$. For the log-transformed CRPS, the optimal value is 0.9, i.e. there is an incentive to report a forecast that is too sharp. **The score is therefore no longer proper.**

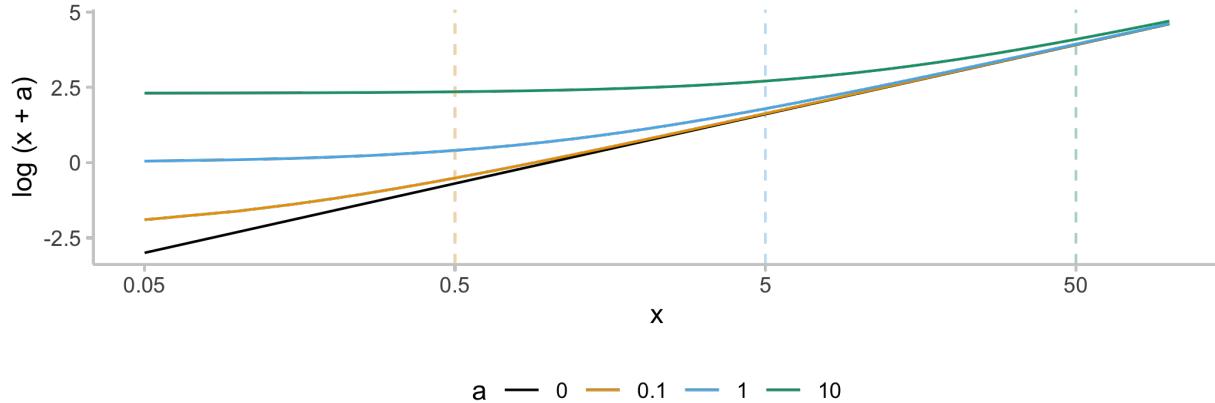


Figure SI.2: Illustration of the effect of adding a small quantity to a value before taking the natural logarithm. For increasing x , all lines eventually approach the black line (representing a transformation with no offset applied). For a given solid line, the dashed line of the same colour marks the x -value that is equal to 5 times the corresponding offset. [It can be seen that for \$a\$ values smaller than one fifth of the transformed quantity, the effect of adding an offset is generally small. When choosing a suitable \$a\$, the trade-off is between staying close to the interpretation of a pure log-transformation \(choosing a small \$a\$ \) and not giving excessive weights to small observations \(by choosing a larger \$a\$, see Figure 6\).](#)

target	<code>type</code>	quantity	measure	natural	log
Cases		Observations	mean	61979	9.19
Cases		Observations	sd	171916	2.10
Cases		Observations	var	29555122130	4.42
Deaths		Observations	mean	220	3.89
Deaths		Observations	sd	435	1.96
Deaths		Observations	var	189051	3.83
Cases		WIS	mean	15840	0.27
Cases		WIS	sd	53117	0.28
Deaths		WIS	mean	31	0.23
Deaths		WIS	sd	65	0.28

Table SI.1: Summary statistics for observations and scores for forecasts from the ECDC data set.

True value	&	Median prediction
> 0		$> 100 \times$ true value
> 10		$> 20 \times$ true value
> 50		$< 1/50 \times$ true value
$= 0$		> 100

Table SI.2: Criteria for removing forecasts. Any forecast that met one of the listed criteria (represented by a row in the table), was removed. Those forecasts were removed in order to be better able to illustrate the effects of the log-transformation on scores and eliminating distortions caused by outlier forecasters. When evaluating models against each other (rather than illustrating the effect of a transformation), one would prefer not to condition on the outcome when deciding whether a forecast should be taken into account.

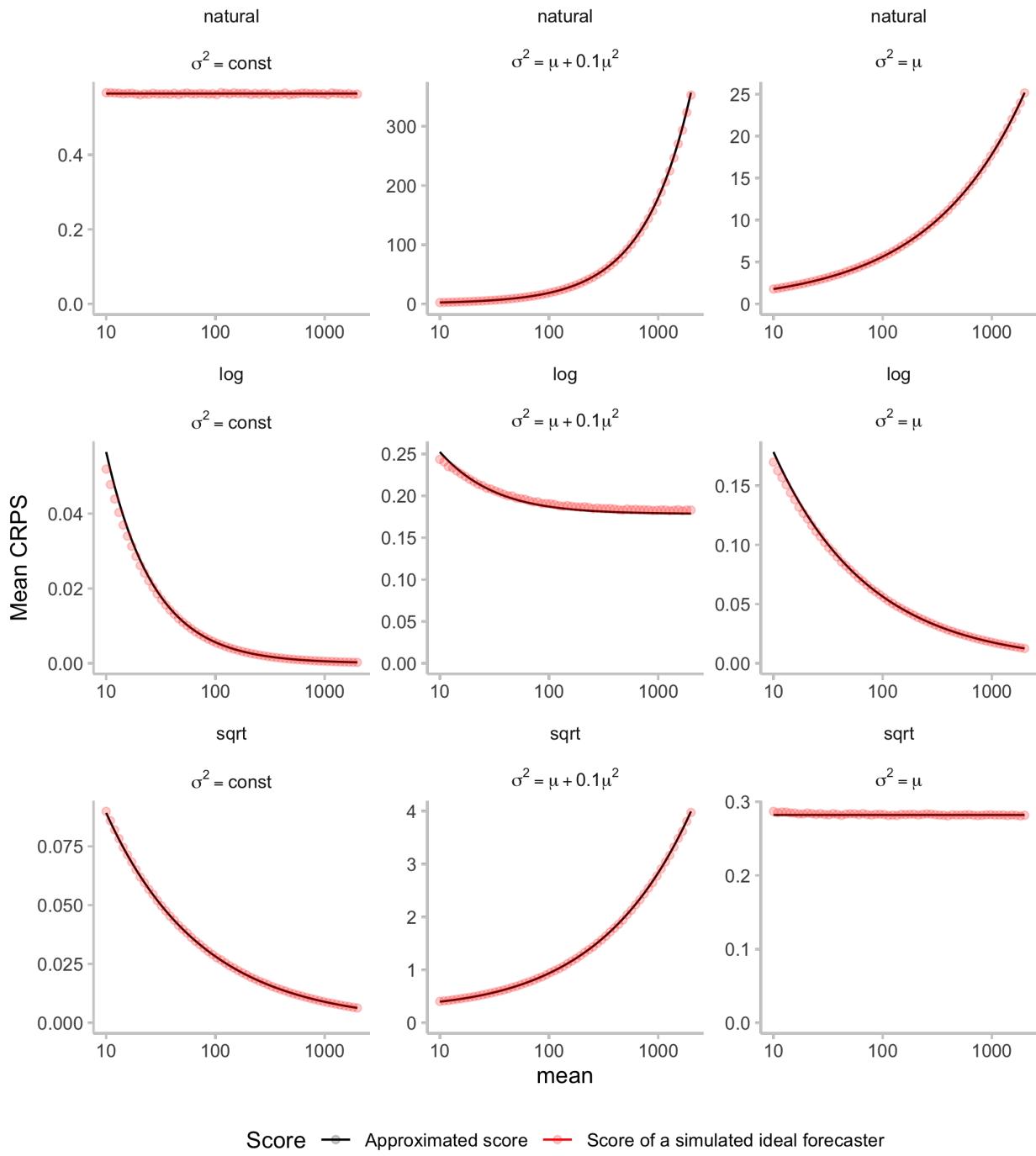


Figure SI.3: Visualisation of expected CRPS values against approximated scores using the approximation detailed in Section 2.4 (see also Figure 2). Expected CRPS scores are shown for three different distributions once on the natural scale (top row) and once scored on the log scale (bottom row).

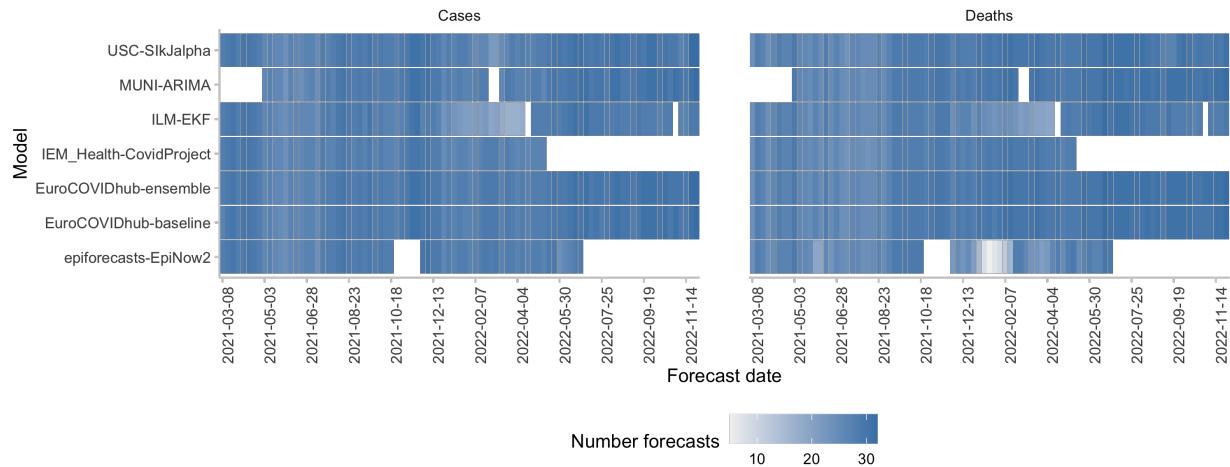


Figure SI.4: Number of forecasts available from different models for each forecast date.

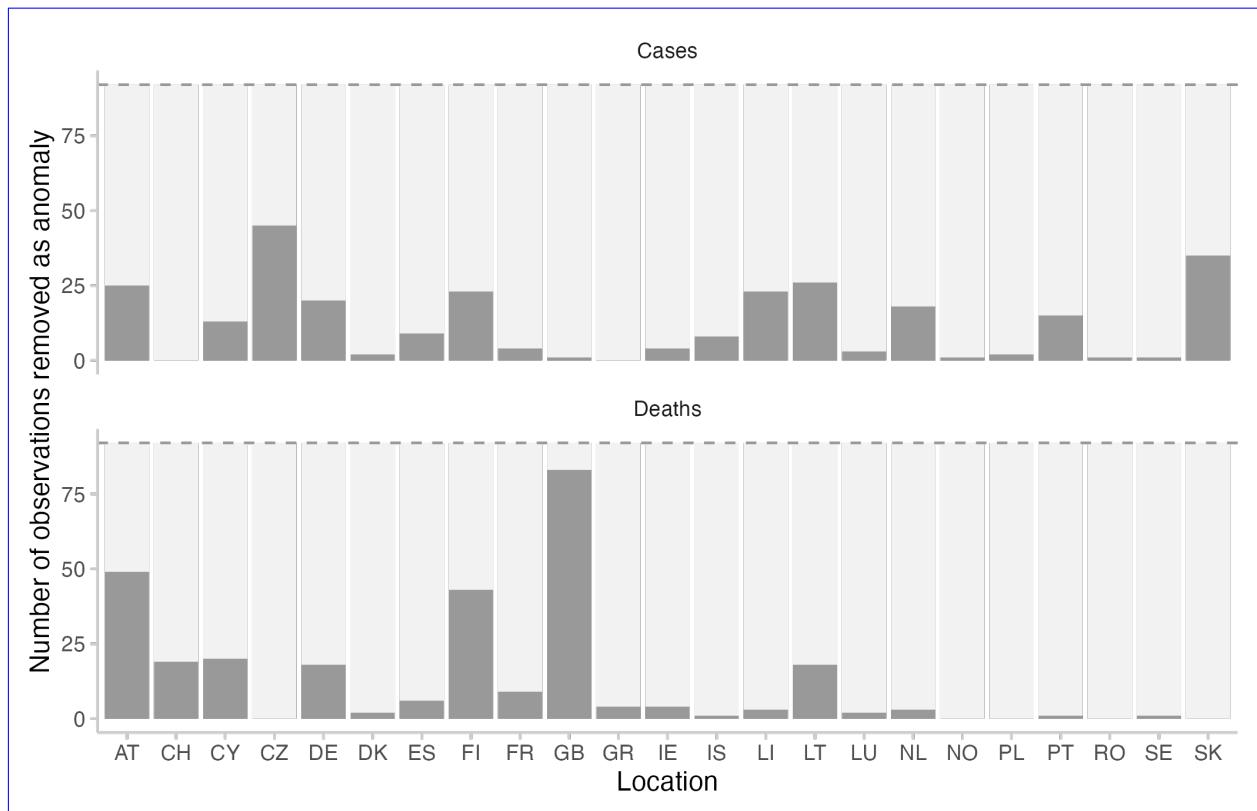


Figure SI.5: Number of observed values that were removed as anomalous. The values were marked as anomalous by the European Forecast Hub team.

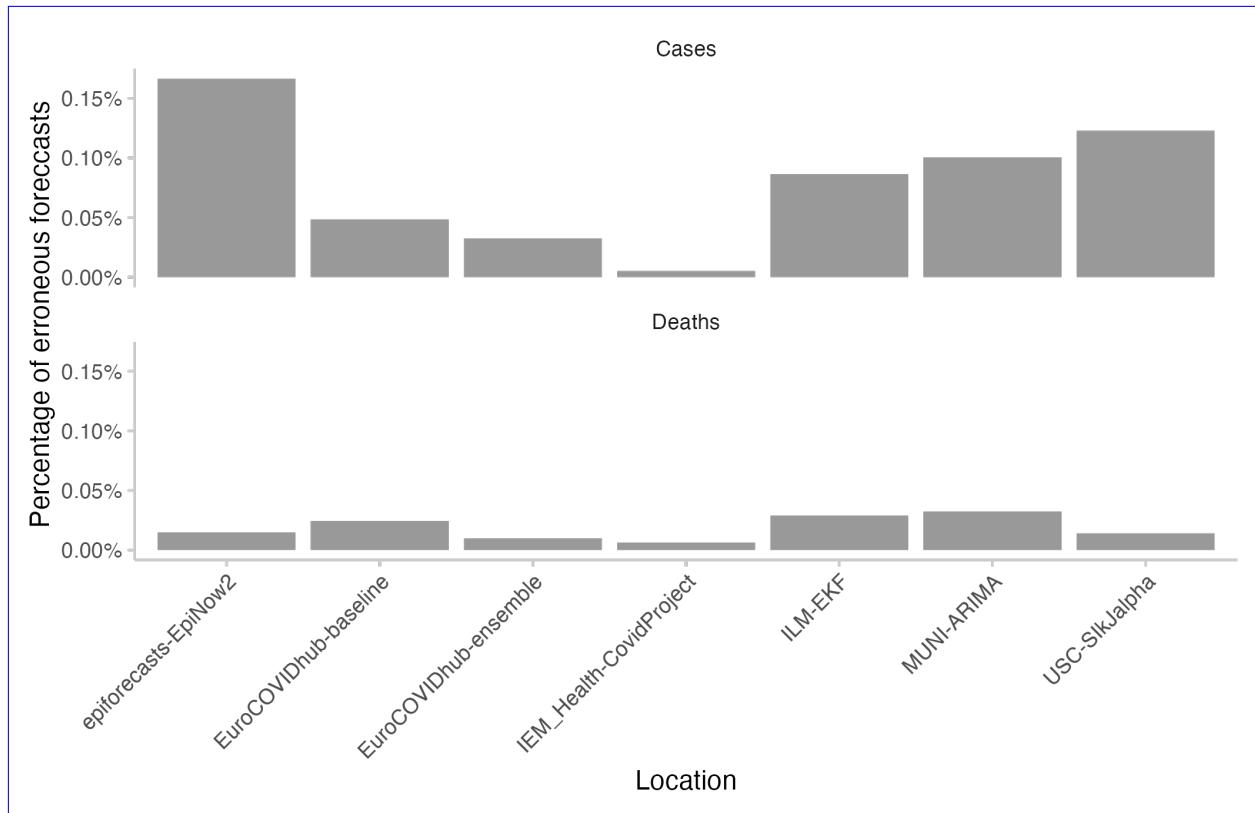


Figure SI.6: Number of forecasts marked as erroneous and removed. Forecasts that were in extremely poor agreement with the observed values were removed from the analysis according to the criteria shown in Table SI.2.

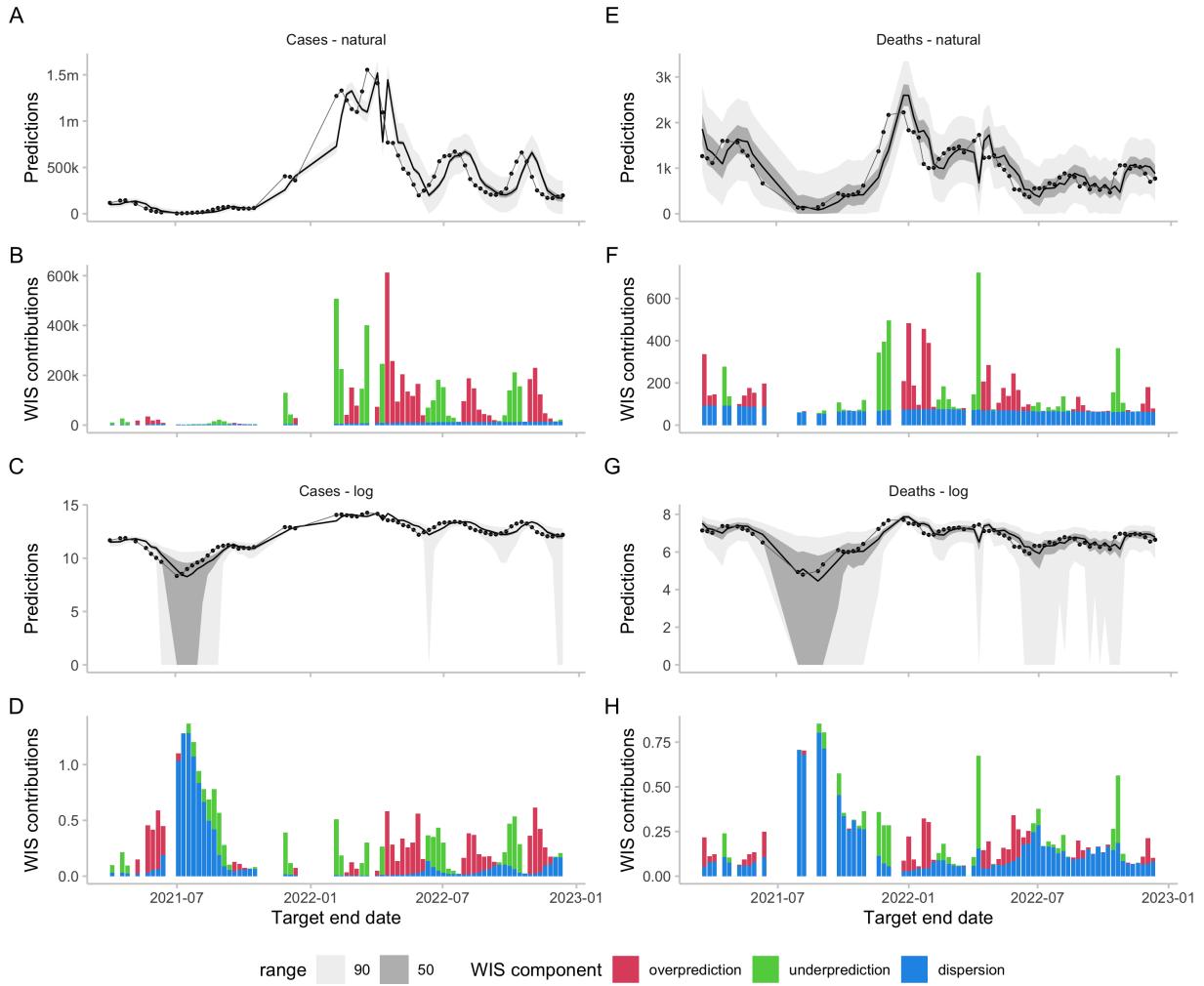


Figure SI.7: Forecasts and scores for two-week-ahead predictions from the EuroCOVIDhub-baseline made in Germany. The model had zero included in some of its 50 percent intervals (e.g. for case forecasts in July 2021), leading to excessive dispersion values on the log scale. One could argue that including zero in the prediction intervals constituted an unreasonable forecast that was rightly penalised, but in general care has to be taken with small numbers. One potential way to do deal with this could be to use a higher a value when applying a transformation $\log(x + a)$, for example $a = 10$ instead of $a = 1$. A, E: 50% and 90% prediction intervals and observed values for cases and deaths on the natural scale. B, F: Corresponding scores. C, G: Forecasts and observations on the log scale. D, H: Corresponding scores.

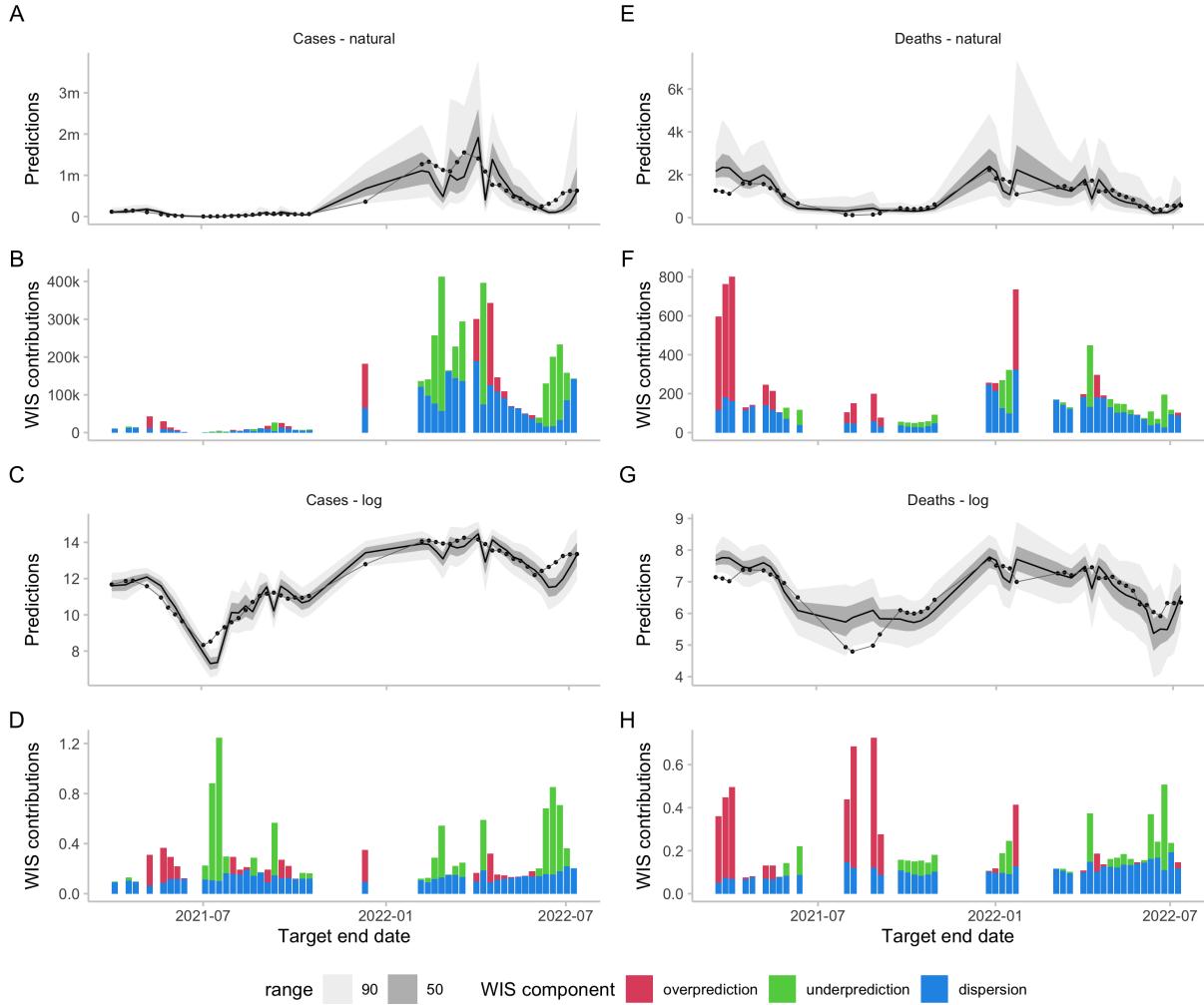


Figure SI.8: Forecasts and scores for two-week-ahead predictions from the epiforecasts-EpiNow2 model (Abbott et al., 2020) made in Germany. A, E: 50% and 90% prediction intervals and observed values for cases and deaths on the natural scale. B, F: Corresponding scores. C, G: Forecasts and observations on the log scale. D, H: Corresponding scores.

Mean WIS in different locations for different transformations applied before scoring. Shown are scores for two-week-ahead forecasts of the EuroCOVIDhub ensemble. On the natural scale (with no transformation prior applying the WIS), scores correlate strongly with the average number of observed values in a given location. The same is true for scores obtained after applying a square-root transformation, or after applying a log-transformation with a large offset a . For illustrative purposes, a was chosen to be 101630 for cases and 530 for deaths, 10 times the respective median observed value. For large values of a , $\log(x + a)$ grows linearly in x , meaning that we expect to observe the same patterns as in the case with no transformation. For decreasing values of a , we give more relative weight to scores in small locations.

608 **References**

- 609 Abbott, S., Hellewell, J., Sherratt, K., Gostic, K., Hickson, J., Badr, H. S., DeWitt, M., Thompson, R.,
610 EpiForecasts, and Funk, S. (2020). *EpiNow2: Estimate Real-Time Case Counts and Time-Varying Epi-*
611 *demiological Parameters*. R package, <https://doi.org/10.5281/zenodo.3957490>.
- 612 Abbott, S., Sherratt, K., Bosse, N., Gruson, H., Bracher, J., and Funk, S. (2022). Evaluating an epidemi-
613 logically motivated surrogate model of a multi-model ensemble.
- 614 Bartlett, M. S. (1936). The Square Root Transformation in Analysis of Variance. *Supplement to the Journal*
615 *of the Royal Statistical Society*, 3(1):68–78.
- 616 Bellégo, C., Benatia, D., and Pape, L. (2022). Dealing with Logs and Zeros in Regression Models.
- 617 Bolin, D. and Wallin, J. (2023). Local scale invariance and robustness of proper scoring rules. *Statistical*
618 *Science*, 38(1):140–159. DOI: 10.1214/22-STS864.
- 619 Bosse, N. I., Gruson, H., Cori, A., van Leeuwen, E., Funk, S., and Abbott, S. (2022). Evaluating Forecasts
620 with scoringutils in R.
- 621 Box, G. E. P. and Cox, D. R. (1964). An Analysis of Transformations. *Journal of the Royal Statistical*
622 *Society. Series B (Methodological)*, 26(2):211–252.
- 623 Bracher, J., Ray, E. L., Gneiting, T., and Reich, N. G. (2021a). Evaluating epidemic forecasts in an interval
624 format. *PLoS computational biology*, 17(2):e1008618.
- 625 Bracher, J., Wolffram, D., Deuschedel, J., Goergen, K., Ketterer, J. L., Ullrich, A., Abbott, S., Barbarossa,
626 M. V., Bertsimas, D., Bhatia, S., Bodych, M., Bosse, N. I., Burgard, J. P., Castro, L., Fairchild, G., Fiedler,
627 J., Fuhrmann, J., Funk, S., Gamin, A., Gogolewski, K., Heyder, S., Hotz, T., Kheifetz, Y., Kirsten, H.,
628 Krueger, T., Krymova, E., Leithäuser, N., Li, M. L., Meinke, J. H., Miasojedow, B., Michaud, I. J.,
629 Mohring, J., Nouvellet, P., Nowosielski, J. M., Ozanski, T., Radwan, M., Rakowski, F., Scholz, M., Soni,
630 S., Srivastava, A., Gneiting, T., and Schienle, M. (2022). National and subnational short-term forecasting
631 of COVID-19 in Germany and Poland, early 2021. *Communications Medicine*. DOI: 10.1038/s43856-022-
632 00191-8.
- 633 Bracher, J., Wolffram, D., Deuschedel, J., Görgen, K., Ketterer, J. L., Ullrich, A., Abbott, S., Barbarossa,
634 M. V., Bertsimas, D., Bhatia, S., Bodych, M., Bosse, N. I., Burgard, J. P., Castro, L., Fairchild, G.,
635 Fuhrmann, J., Funk, S., Gogolewski, K., Gu, Q., Heyder, S., Hotz, T., Kheifetz, Y., Kirsten, H., Krueger,
636 T., Krymova, E., Li, M. L., Meinke, J. H., Michaud, I. J., Niedzielewski, K., Ożański, T., Rakowski,
637 F., Scholz, M., Soni, S., Srivastava, A., Zieliński, J., Zou, D., Gneiting, T., and Schienle, M. (2021b).
638 Short-term forecasting of COVID-19 in Germany and Poland during the second wave – a preregistered
639 study. *medRxiv*, page 2020.12.24.20248826.
- 640 CDC (2022). Cdcepi/Flusight-forecast-data. CDC Epidemic Prediction Initiative. Data repository, <https://github.com/cdcepi/Flusight-forecast-data>.
- 641 Cramer, E., Ray, E. L., Lopez, V. K., Bracher, J., Brennen, A., Rivadeneira, A. J. C., Gerding, A., Gneiting,
642 T., House, K. H., Huang, Y., Jayawardena, D., Kanji, A. H., Khandelwal, A., Le, K., Mühlemann, A.,
643 Niemi, J., Shah, A., Stark, A., Wang, Y., Wattanachit, N., Zorn, M. W., Gu, Y., Jain, S., Bannur, N.,
644 Deva, A., Kulkarni, M., Merugu, S., Raval, A., Shingi, S., Tiwari, A., White, J., Woody, S., Dahan, M.,
645 Fox, S., Gaither, K., Lachmann, M., Meyers, L. A., Scott, J. G., Tec, M., Srivastava, A., George, G. E.,
646 Cegan, J. C., Dettwiller, I. D., England, W. P., Farthing, M. W., Hunter, R. H., Lafferty, B., Linkov,
647 I., Mayo, M. L., Parno, M. D., Rowland, M. A., Trump, B. D., Corsetti, S. M., Baer, T. M., Eisenberg,
648 M. C., Falb, K., Huang, Y., Martin, E. T., McCauley, E., Myers, R. L., Schwarz, T., Sheldon, D., Gibson,
649 G. C., Yu, R., Gao, L., Ma, Y., Wu, D., Yan, X., Jin, X., Wang, Y.-X., Chen, Y., Guo, L., Zhao, Y., Gu,
650 Q., Chen, J., Wang, L., Xu, P., Zhang, W., Zou, D., Biegel, H., Lega, J., Snyder, T. L., Wilson, D. D.,
651 McConnell, S., Walraven, R., Shi, Y., Ban, X., Hong, Q.-J., Kong, S., Turtle, J. A., Ben-Nun, M., Riley,
652 P., Riley, S., Koyluoglu, U., DesRoches, D., Hamory, B., Kyriakides, C., Leis, H., Milliken, J., Moloney,
653

- 654 M., Morgan, J., Ozcan, G., Schrader, C., Shakhnovich, E., Siegel, D., Spatz, R., Stiefeling, C., Wilkinson,
 655 B., Wong, A., Gao, Z., Bian, J., Cao, W., Ferres, J. L., Li, C., Liu, T.-Y., Xie, X., Zhang, S., Zheng,
 656 S., Vespignani, A., Chinazzi, M., Davis, J. T., Mu, K., y Piontti, A. P., Xiong, X., Zheng, A., Baek, J.,
 657 Farias, V., Georgescu, A., Levi, R., Sinha, D., Wilde, J., Penna, N. D., Celi, L. A., Sundar, S., Cavany,
 658 S., España, G., Moore, S., Oidtmann, R., Perkins, A., Osthus, D., Castro, L., Fairchild, G., Michaud, I.,
 659 Karlen, D., Lee, E. C., Dent, J., Grantz, K. H., Kaminsky, J., Kaminsky, K., Keegan, L. T., Lauer, S. A.,
 660 Lemaitre, J. C., Lessler, J., Meredith, H. R., Perez-Saez, J., Shah, S., Smith, C. P., Truelove, S. A., Wills,
 661 J., Kinsey, M., Obrecht, R. F., Tallaksen, K., Burant, J. C., Wang, L., Gao, L., Gu, Z., Kim, M., Li,
 662 X., Wang, G., Wang, Y., Yu, S., Reiner, R. C., Barber, R., Gaikedu, E., Hay, S., Lim, S., Murray, C.,
 663 Pigott, D., Prakash, B. A., Adhikari, B., Cui, J., Rodríguez, A., Tabassum, A., Xie, J., Keskinocak, P.,
 664 Asplund, J., Baxter, A., Oruc, B. E., Serban, N., Arik, S. O., Dusenberry, M., Epshteyn, A., Kanal, E.,
 665 Le, L. T., Li, C.-L., Pfister, T., Sava, D., Sinha, R., Tsai, T., Yoder, N., Yoon, J., Zhang, L., Abbott,
 666 S., Bosse, N. I., Funk, S., Hellewel, J., Meakin, S. R., Munday, J. D., Sherratt, K., Zhou, M., Kalantari,
 667 R., Yamana, T. K., Pei, S., Shaman, J., Ayer, T., Adee, M., Chhatwal, J., Dalgic, O. O., Ladd, M. A.,
 668 Linas, B. P., Mueller, P., Xiao, J., Li, M. L., Bertsimas, D., Lami, O. S., Soni, S., Bouardi, H. T., Wang,
 669 Y., Wang, Q., Xie, S., Zeng, D., Green, A., Bien, J., Hu, A. J., Jahja, M., Narasimhan, B., Rajanala, S.,
 670 Rumack, A., Simon, N., Tibshirani, R., Tibshirani, R., Ventura, V., Wasserman, L., O'Dea, E. B., Drake,
 671 J. M., Pagano, R., Walker, J. W., Slayton, R. B., Johansson, M., Biggerstaff, M., and Reich, N. G. (2021).
 672 Evaluation of individual and ensemble probabilistic forecasts of COVID-19 mortality in the US. *medRxiv*,
 673 page 2021.02.03.21250974.
- 674 Cramer, E., Reich, N. G., Wang, S. Y., Niemi, J., Hannan, A., House, K., Gu, Y., Xie, S., Horstman,
 675 S., aniruddhadiga, Walraven, R., starkari, Li, M. L., Gibson, G., Castro, L., Karlen, D., Wattanachit,
 676 N., jinghuichen, zyt9lsb, agarwal1996, Woody, S., Ray, E., Xu, F. T., Biegel, H., GuidoEspana, X., X.,
 677 Bracher, J., Lee, E., har96, and leyouz (2020). COVID-19 Forecast Hub: 4 December 2020 snapshot.
- 678 Diks, C., Panchenko, V., and van Dijk, D. (2011). Likelihood-based scoring rules for comparing density
 679 forecasts in tails. *Journal of Econometrics*, 163(2):215–230.
- 680 Dunn, P. K. and Smyth, G. K. (2018). *Generalized Linear Models With Examples in R*. Springer.
- 681 Dushoff, J. and Park, S. W. (2021). Speed and strength of an epidemic intervention. *Proceedings of the
 682 Royal Society B: Biological Sciences*, 288(1947):20201556.
- 683 European Covid-19 Forecast Hub (2021). European Covid-19 Forecast Hub. <https://covid19forecasthub.eu/>.
- 684 Flores, B. E. (1986). A pragmatic view of accuracy measurement in forecasting. *Omega*, 14(2):93–98.
- 685 Fuglstad, G.-A., Simpson, D., Lindgren, F., and Rue, H. (2015). Does non-stationary spatial data always
 686 require non-stationary random fields? *Spatial Statistics*, 14:505–531.
- 687 Funk, S., Camacho, A. J., Kucharski, A. J., Lowe, R., Eggo, R. M., and Edmunds, W. J. (2019). Assessing the
 688 performance of real-time epidemic forecasts: A case study of Ebola in the Western Area region of Sierra
 689 Leone, 2014–15. *PLOS Computational Biology*, 15(2):e1006785.
- 690 Gneiting, T. (2011). Making and Evaluating Point Forecasts. *Journal of the American Statistical Association*,
 691 106(494):746–762.
- 692 Gneiting, T. and Raftery, A. E. (2005). Weather Forecasting with Ensemble Methods. *Science*,
 693 310(5746):248–249.
- 694 Gneiting, T. and Raftery, A. E. (2007). Strictly Proper Scoring Rules, Prediction, and Estimation. *Journal
 695 of the American Statistical Association*, 102(477):359–378.
- 696 Good, I. J. (1952). Rational Decisions. *Journal of the Royal Statistical Society. Series B (Methodological)*,
 697 14(1):107–114.

- 698 Gostic, K. M., McGough, L., Baskerville, E., Abbott, S., Joshi, K., Tedijanto, C., Kahn, R., Niehus, R.,
699 Hay, J., de Salazar, P., Hellewell, J., Meakin, S., Munday, J., Bosse, N. I., Sherrat, K., Thompson, R. N.,
700 White, L. F., Huisman, J. S., Scire, J., Bonhoeffer, S., Stadler, T., Wallinga, J., Funk, S., Lipsitch, M., and
701 Cobey, S. (2020). Practical considerations for measuring the effective reproductive number, *Rt*. *medRxiv*.
- 702 Guerrero, V. M. (1993). Time-series analysis supported by power transformations. *Journal of Forecasting*,
703 12(1):37–48.
- 704 Held, L., Meyer, S., and Bracher, J. (2017). Probabilistic forecasting in infectious disease epidemiology: The
705 13th Armitage lecture. *Statistics in Medicine*, 36(22):3443–3460.
- 706 Johansson, M. A., Apfeldorf, K. M., Dobson, S., Devita, J., Buczak, A. L., Baugher, B., Moniz, L. J.,
707 Bagley, T., Babin, S. M., Guven, E., Yamana, T. K., Shaman, J., Moschou, T., Lothian, N., Lane, A.,
708 Osborne, G., Jiang, G., Brooks, L. C., Farrow, D. C., Hyun, S., Tibshirani, R. J., Rosenfeld, R., Lessler,
709 J., Reich, N. G., Cummings, D. A. T., Lauer, S. A., Moore, S. M., Clapham, H. E., Lowe, R., Bailey,
710 T. C., García-Díez, M., Carvalho, M. S., Rodó, X., Sardar, T., Paul, R., Ray, E. L., Sakrejda, K., Brown,
711 A. C., Meng, X., Osoba, O., Vardavas, R., Manheim, D., Moore, M., Rao, D. M., Porco, T. C., Ackley,
712 S., Liu, F., Worden, L., Convertino, M., Liu, Y., Reddy, A., Ortiz, E., Rivero, J., Brito, H., Juarrero, A.,
713 Johnson, L. R., Gramacy, R. B., Cohen, J. M., Mordecai, E. A., Murdock, C. C., Rohr, J. R., Ryan, S. J.,
714 Stewart-Ibarra, A. M., Weikel, D. P., Jutla, A., Khan, R., Poultney, M., Colwell, R. R., Rivera-García,
715 B., Barker, C. M., Bell, J. E., Biggerstaff, M., Swerdlow, D., Mier-y-Teran-Romero, L., Forshey, B. M.,
716 Trtanj, J., Asher, J., Clay, M., Margolis, H. S., Hebbeler, A. M., George, D., and Jean-Paul Chretien
717 (2019). An open challenge to advance probabilistic forecasting for dengue epidemics. *Proceedings of the
718 National Academy of Sciences*, 116(48):24268–24274.
- 719 Lehmann, E. L. (1950). Some Principles of the Theory of Testing Hypotheses. *The Annals of Mathematical
720 Statistics*, 21(1):1 – 26.
- 721 Lerch, S., Thorarinsdottir, T. L., Ravazzolo, F., and Gneiting, T. (2015). Forecaster’s Dilemma: Extreme
722 Events and Forecast Evaluation.
- 723 Löwe, R., Mikkelsen, P. S., and Madsen, H. (2014). Stochastic rainfall-runoff forecasting: Parameter estimation,
724 multi-step prediction, and evaluation of overflow risk. *Stochastic Environmental Research and Risk
725 Assessment*, 28(3):505–516.
- 726 Mayr, J. and Ulbricht, D. (2015). Log versus level in VAR forecasting: 42 million empirical answers—Expect
727 the unexpected. *Economics Letters*, 126:40–42.
- 728 Pellis, L., Scarabel, F., Stage, H. B., Overton, C. E., Chappell, L. H. K., Fearon, E., Bennett, E., Lythgoe,
729 K. A., House, T. A., Hall, I., and null, n. (2021). Challenges in control of COVID-19: Short doubling time
730 and long delay to effect of interventions. *Philosophical Transactions of the Royal Society B: Biological
731 Sciences*, 376(1829):20200264.
- 732 R Core Team (2022). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical
733 Computing, Vienna, Austria.
- 734 Reich, N. G., Brooks, L. C., Fox, S. J., Kandula, S., McGowan, C. J., Moore, E., Osthus, D., Ray, E. L.,
735 Tushar, A., Yamana, T. K., Biggerstaff, M., Johansson, M. A., Rosenfeld, R., and Shaman, J. (2019).
736 A collaborative multiyear, multimodel assessment of seasonal influenza forecasting in the United States.
737 *Proceedings of the National Academy of Sciences*, 116(8):3146–3154.
- 738 Reich, N. G., Lessler, J., Funk, S., Viboud, C., Vespignani, A., Tibshirani, R. J., Shea, K., Schienle, M.,
739 Runge, M. C., Rosenfeld, R., Ray, E. L., Niehus, R., Johnson, H. C., Johansson, M. A., Hochheiser, H.,
740 Gardner, L., Bracher, J., Borcherding, R. K., and Biggerstaff, M. (2022). Collaborative hubs: Making the
741 most of predictive epidemic modeling. *American Journal of Public Health*, 112(6):839–842.

- 742 Sherratt, K., Gruson, H., Grah, R., Johnson, H., Niehus, R., Prasse, B., Sandman, F., Deusel, J., Wolffram,
 743 D., Abbott, S., Ullrich, A., Gibson, G., Ray, EL., Reich, NG., Sheldon, D., Wang, Y., Wattanachit, N.,
 744 Wang, L., Trnka, J., Obozinski, G., Sun, T., Thanou, D., Pottier, L., Krymova, E., Barbarossa, MV.,
 745 Leithäuser, N., Mohring, J., Schneider, J., Wlazlo, J., Fuhrmann, J., Lange, B., Rodiah, I., Baccam,
 746 P., Gurung, H., Stage, S., Suchoski, B., Budzinski, J., Walraven, R., Villanueva, I., Tucek, V., Šmíd, M.,
 747 Zajícek, M., Pérez, Á. C., Reina, B., Bosse, NI., Meakin, S., Di Loro, A., Maruotti, A., Eclerová, V., Kraus,
 748 A., Kraus, D., Pribylova, L., Dimitris, B., Li, ML., Saksham, S., Dehning, J., Mohr, S., Priesemann, V.,
 749 Redlarski, G., Bejar, B., Ardenghi, G., Parolini, N., Ziarelli, G., Bock, W., Heyder, S., Hotz, T., E., S. D.,
 750 Guzman-Merino, M., Aznarte, JL., Moriña, D., Alonso, S., Alvarez, E., López, D., Prats, C., Burgard, JP.,
 751 Rodloff, A., Zimmermann, T., Kuhlmann, A., Zibert, J., Pennoni, F., Divino, F., Català, M., Lovison,
 752 G., Giudici, P., Tarantino, B., Bartolucci, F., Jona, L. G., Mingione, M., Farcomeni, A., Srivastava,
 753 A., Montero-Manso, P., Adiga, A., Hurt, B., Lewis, B., Marathe, M., Porebski, P., Venkatraman, S.,
 754 Bartczuk, R., Dreger, F., Gambin, A., Gogolewski, K., Gruziel-Słomka, M., Krupa, B., Moszynski, A.,
 755 Niedzielewski, K., Nowosielski, J., Radwan, M., Rakowski, F., Semeniuk, M., Szczurek, E., Zielinski, J.,
 756 Kisielewski, J., Pabjan, B., Holger, K., Kheifetz, Y., Scholz, M., Bodych, M., Filinski, M., Idzikowski,
 757 R., Krueger, T., Ozanski, T., Bracher, J., and Funk, S. (2022). Predictive performance of multi-model
 758 ensemble forecasts of COVID-19 across European nation.
- 759 Srivastava, A., Singh, S., and Lee, F. (2022). Shape-based Evaluation of Epidemic Forecasts.
- 760 Taylor, J. W. (1999). Evaluating volatility and interval forecasts. *Journal of Forecasting*, 18(2):111–128.
- 761 Timmermann, A. (2018). Forecasting Methods in Finance. *Annual Review of Financial Economics*,
 762 10(1):449–479.
- 763 Wallinga, J. and Lipsitch, M. (2007). How generation intervals shape the relationship between growth rates
 764 and reproductive numbers. *Proceedings of the Royal Society B: Biological Sciences*, 274(1609):599–604.
- 765 Winkler, R. (1996). Scoring rules and the evaluation of probabilities. *TEST*, 5(1):1–60.