

Comparing trained and untrained probabilistic ensemble forecasts of COVID-19 cases and deaths in the United States

Evan L. Ray^{a,*}, Logan C. Brooks^b, Jacob Bien^c, Matthew Biggerstaff^d, Nikos I. Bosse^e, Johannes Bracher^{f,g}, Estee Y. Cramer^a, Sebastian Funk^e, Aaron Gerding^a, Michael A. Johansson^d, Aaron Rumack^b, Yijin Wang^a, Martha Zorn^a, Ryan J. Tibshirani^b, Nicholas G. Reich^a

^a*School of Public Health and Health Sciences, University of Massachusetts Amherst*

^b*Machine Learning Department, Carnegie Mellon University*

^c*Department of Data Sciences and Operations, University of Southern California*

^d*COVID-19 Response, U.S. Centers for Disease Control and Prevention*

^e*London School of Hygiene & Tropical Medicine*

^f*Chair of Statistical Methods and Econometrics, Karlsruhe Institute of Technology*

^g*Computational Statistics Group, Heidelberg Institute for Theoretical Studies*

Abstract

The U.S. COVID-19 Forecast Hub aggregates forecasts of the short-term burden of COVID-19 in the United States from many contributing teams. We study methods for building an ensemble that combines forecasts from these teams. These experiments have informed the ensemble methods used by the Hub. To be most useful to policy makers, ensemble forecasts must have stable performance in the presence of two key characteristics of the component forecasts: (1) occasional misalignment with the reported data, and (2) instability in the relative performance of component forecasters over time. Our results indicate that in the presence of these challenges, an untrained and robust approach to ensembling using an equally weighted median of all component forecasts is a good choice to support public health decision makers. In settings where some contributing forecasters have a stable record of good performance, trained ensembles that give those forecasters higher weight can also be helpful.

Keywords: Health forecasting, epidemiology, COVID-19, ensemble, quantile combination

*Corresponding author

Email address: e1ray@umass.edu (Evan L. Ray)

Comparing trained and untrained probabilistic ensemble forecasts of COVID-19 cases and deaths in the United States

Abstract

The U.S. COVID-19 Forecast Hub aggregates forecasts of the short-term burden of COVID-19 in the United States from many contributing teams. We study methods for building an ensemble that combines forecasts from these teams. These experiments have informed the ensemble methods used by the Hub. To be most useful to policy makers, ensemble forecasts must have stable performance in the presence of two key characteristics of the component forecasts: (1) occasional misalignment with the reported data, and (2) instability in the relative performance of component forecasters over time. Our results indicate that in the presence of these challenges, an untrained and robust approach to ensembling using an equally weighted median of all component forecasts is a good choice to support public health decision makers. In settings where some contributing forecasters have a stable record of good performance, trained ensembles that give those forecasters higher weight can also be helpful.

Keywords: Health forecasting, epidemiology, COVID-19, ensemble, quantile combination

1. Introduction

Accurate short-term forecasts of infectious disease indicators (i.e., disease surveillance signals) can inform public health decision-making and outbreak response activities such as non-pharmaceutical interventions, site selection for clinical trials of pharmaceutical treatments, or the distribution of limited health care resources (Wallinga et al., 2010; Lipsitch et al., 2011; Dean et al., 2020). Epidemic forecasts have been incorporated into public health decision-making in a wide variety of situations, including outbreaks of dengue fever in Brazil, Vietnam, and Thailand (Lowe et al., 2016; Coln-Gonzalez et al., 2021; Reich et al., 2016) and influenza in the U.S. (McGowan et al., 2019).

These efforts frequently use ensemble forecasts that combine predictions from many models. In a wide array of fields, ensemble approaches have provided consistent improvements in accuracy

1
2
3 and robustness relative to stand-alone forecasts (Gneiting & Raftery, 2005; Polikar, 2006). The
4 usefulness of ensemble forecasts has also been demonstrated repeatedly in multiple infectious disease
5 settings, including influenza, Ebola, dengue, RSV, and others (Yamana et al., 2016; Viboud et al.,
6 2018; McGowan et al., 2019; Johansson et al., 2019; Reis et al., 2019; Reich et al., 2019). In light of
7 this record of strong performance, ensembles are natural candidates for forecasts used as an input
8 to high-stakes public health decision-making processes.
9

10
11 This paper describes ensemble modeling efforts at the U.S. COVID-19 Forecast Hub (<https://covid19forecasthub.org/>, hereafter the “U.S. Hub”), from spring 2020 through spring 2022.
12 Starting in April 2020, the U.S. Hub created ensemble forecasts of reported incident deaths one
13 through four weeks ahead in the 50 states, Washington, D.C., and 6 territories as well as at the
14 national level by combining forecasts submitted by a large and variable number of contributing
15 teams using different modeling techniques and data sources. In July 2020, forecasts of incident
16 reported COVID-19 cases were added. Of note, the U.S. Hub produces *probabilistic* forecasts
17 in which uncertainty about future disease incidence is quantified through the specification of a
18 predictive distribution that is represented by a collection of predictive quantiles. Since the inception
19 of the U.S. Hub, these ensemble forecasts have been provided to the U.S. Centers for Disease Control
20 and Prevention (CDC) and have been the basis of official CDC forecasting communications (CDC,
21 2021).
22
23
24
25
26
27
28
29
30
31
32
33
34
35

36 1.1. Related literature

37
38 A wide variety of standalone methodological approaches have been shown to be able to make
39 forecasts of short-term outbreak activity that are more accurate than naive baseline forecasts in
40 various epidemiological settings. Some approaches have used existing statistical frameworks to
41 model associations between outcomes of interest and known or hypothesized drivers of outbreaks,
42 such as recent trends in transmission or environmental factors. To cite just a few examples,
43 methods used include multiscale probabilistic Bayesian random walk models (Osthus & Moran,
44 2021), Gaussian processes (Johnson et al., 2018), kernel conditional density estimation (Ray et al.,
45 2017; Brooks et al., 2018), and generalized additive models (Lauer et al., 2018). Other models
46 have an implicit or explicit representation of a disease transmission process, such as variations
47 on the susceptible-infectious-recovered (SIR) compartmental model (Shaman & Karspeck, 2012;
48 Lega & Brown, 2016; Osthus et al., 2017; Pei et al., 2018; Turtle et al., 2021). Aspects of these
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3 modeling frameworks can also be combined, for instance using time series methods to build models
4 that have a compartmental structure or incorporate key epidemiological parameters such as the
5 effective reproduction number R_t , or models that use a time series process to capture systematic
6 deviations from a compartmental core (Bartolucci et al., 2021; Agosto et al., 2021; Osthus et al.,
7 2019).

8
9
10
11
12 There is a large literature on ensemble forecasting, but of particular relevance to the present
13 work is the research on combining, calibrating and evaluating distributional forecasts (Gneiting &
14 Raftery, 2007; Gneiting et al., 2007; Ranjan & Gneiting, 2010; Claeskens et al., 2016). We note
15 that prior work on forecast combination has mostly focused on combining forecasts represented
16 as probability densities or probability mass functions rather than forecasts parameterized by a set
17 of discrete quantile levels, which is the format of the forecasts in the present study. However, in
18 psychological studies there is a long history of combining quantiles from multiple distributions as
19 a mechanism for summarizing distributions of response times, error rates, and similar quantities
20 across many subjects (Vincent, 1912; Ratcliff, 1979). More recently, this approach has also been
21 used to combine probabilistic assessments from multiple subject matter experts or statistical models
22 in fields such as security threat detection and economic forecasting (Hora et al., 2013; Lichtendahl Jr
23 et al., 2013; Gaba et al., 2017; Buseti, 2017). In the context of infectious disease forecasting,
24 Bracher et al. (2021b) conducted a similar but less extensive analysis to the one presented here using
25 data from a related forecast hub focusing on Germany and Poland. Taylor & Taylor (2021) recently
26 explored several approaches to constructing quantile-based ensemble forecasts of cumulative deaths
27 due to COVID-19 using the data from the U.S. Hub, although they did not generate ensemble
28 forecasts in real time or appear to have used the specific versions of ground truth data that were
29 available for constructing ensembles in real time.

30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

As was mentioned earlier, ensemble forecasts have also been used in a variety of other appli-
cations in real-time forecasting of infectious diseases, often with seasonal transmission dynamics
where many years of training data are available (Yamana et al., 2016; Reich et al., 2019; Reis et al.,
2019; Coln-Gonzlez et al., 2021). In such applications, simple combination approaches have gener-
ally been favored over complex ones, with equal-weighted approaches often performing similarly
to trained approaches that assign weights to different models based on past performance (Ray &
Reich, 2018; Bracher et al., 2021b). These results align with theory suggesting that the uncertainty

1
2
3 in weight estimation can pose a challenge in applications with a low signal-to-noise ratio (Claeskens
4 et al., 2016).
5
6

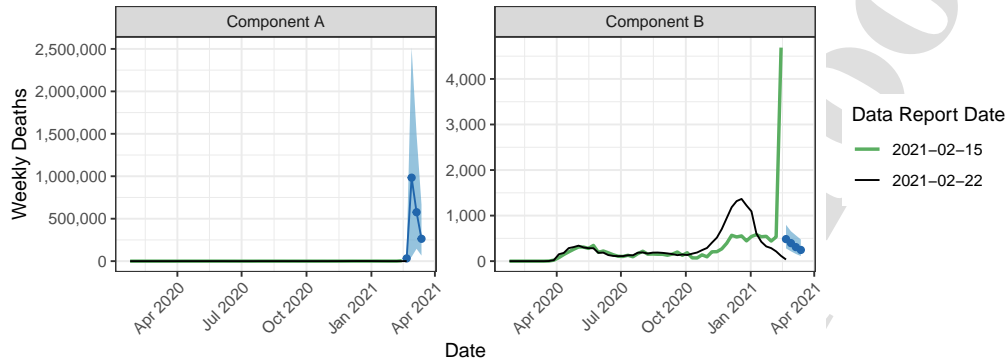
7 8 *1.2. Contributions of this article* 9

10 This paper is focused on explaining the careful considerations that have gone into building
11 a relatively simple “production” ensemble model for a difficult, high-stakes, real-time prediction
12 problem: forecasting COVID-19 cases and deaths in the United States, to support public health
13 decision-making. We do not empirically investigate the performance of complex forecast combi-
14 nation strategies from the online prediction literature, which generally require richer and larger
15 training data sets.
16
17

18 The goal of the U.S. Hub in developing an operational ensemble was to produce forecasts of the
19 short-term trajectory of COVID-19 that had good performance on average and stable performance
20 across time and different locations. Real-time forecasting for an emerging pathogen in an open,
21 collaborative setting introduces important challenges that an ensemble combination method must
22 be able to handle. First, teams occasionally submitted outlying component forecasts due to software
23 errors, incorrect model assumptions, or a lack of robustness to input data anomalies (Figure 1 (a),
24 Supplemental Figures 1 and 2). Second, some component models were generally better than others,
25 but the relative performance of different models was somewhat unstable across time (Figure 1
26 (b), Supplemental Figures 3 and 4). In particular, some forecasters alternated between being
27 among the best-performing models and among the worst-performing models within a span of a few
28 weeks, which introduced a challenge for ensemble methods that attempted to weight component
29 forecasters based on their past performance. In this manuscript, we explore and compare variations
30 on ensemble methods designed to address these challenges and produce real-time forecasts that are
31 as accurate as possible to support public health decision-makers.
32
33

34 We give detailed results from experiments that were run concurrently with the weekly releases of
35 ensemble forecasts from the start of the U.S. Hub in 2020 through the spring of 2022, as documented
36 in preliminary reports (Brooks et al., 2020 [Online]; Ray et al., 2021 [Online]). These experiments
37 provided the evidence for decisions (a) to move to a median-based ensemble from one based on
38 means in July 2020; (b) to switch to a trained ensemble for forecasts of deaths in November 2021;
39 and (c) to implement a weight regularization strategy for that trained ensemble starting in January
40 2022. In a secondary analysis, we also consider the prospective performance of these methods in
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

(a) Forecasts of incident deaths in Ohio from February 15, 2021



(b) Component forecaster relative WIS for forecasts of incident cases in the US

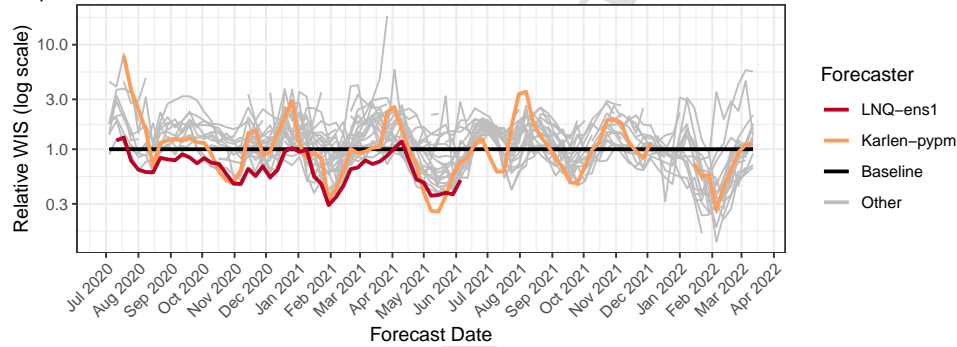


Figure 1: (a) Predictive medians and 95% prediction intervals for incident deaths in Ohio generated on February 15, 2021 by two example component forecasters. The vertical axis scale is different in each facet, reflecting differences across several orders of magnitude in forecasts from different forecasters; the reference data are the same in each plot. The data that were available as of Monday, February 15, 2021 included a large spike in reported deaths that had been redistributed into the history of the time series in the version of the data available as of Monday, February 22, 2021. In this panel, forecaster names are anonymized to avoid calling undue attention to individual teams; similar behavior has been exhibited by many forecasters. (b) Illustration of the relative weighted interval score (WIS, defined in Section 2.5) of component forecasters over time; lower scores indicate better performance. Each point summarizes the skill of forecasts made at a given date for the one through four week ahead forecasts of incident cases across all state-level locations.

1
2
3 the closely related setting of forecasting cases and deaths in Europe, to examine the generalizability
4 of the results from our experiments using data from the U.S.
5
6

7 The following sections document the format and general characteristics of COVID-19 forecasts
8 under consideration, the ensemble approaches studied, and the results of comparing different ap-
9 proaches both during model development and during a prospective evaluation of selected methods.
10
11

12 **2. Methods**

13
14 We give an overview of the U.S. and European Forecast Hubs and the high-level structure of
15 our experiments in Sections 2.1 through 2.5, and then describe the ensemble methods that we
16 consider in Section 2.6.
17
18

19 *2.1. Problem context: forecasting short-term COVID-19 burden*

20
21 Starting in April 2020, the U.S. Hub collected probabilistic forecasts of the short-term burden
22 of COVID-19 in the U.S. at the national, state/territory, and county levels (Cramer et al., 2021a);
23 a similar effort began in February 2021 for forecasts of disease burden in 32 European countries
24 (European Covid-19 Forecast Hub, 2021). In this manuscript, we focus on constructing probabilis-
25 tic ensemble forecasts of weekly counts of reported cases and deaths due to COVID-19 at forecast
26 horizons of one to four weeks for states and territories in the U.S. and for countries in Europe.
27 A maximum horizon of four weeks was set by collaborators at CDC as a horizon at which fore-
28 casts would be useful to public health practitioners while maintaining reasonable expectations of a
29 minimum standard of forecast accuracy and reliability. Probabilistic forecasts were contributed to
30 the Hubs in a quantile-based format by teams in academia, government, and industry. The Hubs
31 produced ensemble forecasts each week on Monday using forecasts from teams contributing that
32 week. In the U.S. Hub, seven quantile levels were used for forecasts of cases and 23 quantile levels
33 were used for forecasts of deaths; in the European Hub, 23 quantile levels were used for both target
34 variables.
35
36

37
38 Weekly reported cases and deaths were calculated as the difference in cumulative counts on
39 consecutive Saturdays, using data assembled by the Johns Hopkins University Center for Systems
40 Science and Engineering as the ground truth (Dong et al., 2020). Due to changes in the definitions
41 of reportable cases and deaths, as well as errors in reporting and backlogs in data collection, there
42 were some instances in which the ground truth data included outlying values, or were revised. Most
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3 outliers and revisions were inconsequential, but some were quite substantial in the U.S. as well as in
4 Europe (Figure 2). When fitting retrospective ensembles, we fit to the data that would have been
5 available in real time. This is critical because the relative performance of different component fore-
6 casters may shift dramatically depending on whether originally-reported or subsequently-revised
7 data were used to measure forecast skill. An ensemble trained using revised data can therefore
8 have a substantial advantage over one trained using only data that were available in real time, and
9 its performance is not a reliable gauge of how that ensemble method might have done in real time.

10
11
12
13
14
15
16 The U.S. Hub conducted extensive ensemble model development in real time from late July
17 2020 through the end of April 2021, with smaller focused experiments ongoing thereafter. We
18 present results for the model development phase as well as a prospective evaluation of a subset
19 of ensemble methods in the U.S. starting with forecasts created on May 3, 2021 and continuing
20 through March 14, 2022. We note that we continued examining a wider range of methods to inform
21 weekly operational forecasting tasks, but the methods that we chose to evaluate prospectively were
22 selected by May 3, 2021, the beginning of the prospective evaluation period, with no alterations
23 thereafter. Real-time submissions of the relative WIS weighted median ensemble described below
24 are on record in the U.S. Hub for the duration of the prospective evaluation period. In one
25 section of the results below, we present a small post hoc exploration of the effects of regularizing
26 the component forecaster weights; these results should be interpreted with caution as they do not
27 constitute a prospective evaluation. To examine how well our findings generalize, we also evaluated
28 the performance of a subset of ensemble methods for prospective forecasts of cases and deaths at
29 the national level for countries in Europe from May 3, 2021 to March 14, 2022.

30 31 32 33 34 35 36 37 38 39 40 41 42 *2.2. Eligibility criteria*

43
44 In the Forecast Hubs, not all forecasts from contributing models are available for all weeks. For
45 example, forecasters may have started submitting forecasts in different weeks, and some forecasters
46 submitted forecasts for only a subset of locations in one or more weeks.

47
48
49 The ensemble forecast for a particular location and forecast date included all component fore-
50 casts with a complete set of predictive quantiles (i.e., 7 predictive quantiles for incident cases, 23
51 for deaths) for all 4 forecast horizons. Teams were not required to submit forecasts for all locations
52 to be included in the ensemble. Some ensemble methods that we considered require historical
53 forecasts to inform component model selection or weighting; for these methods, at least one prior
54
55
56
57
58
59
60
61
62
63
64
65

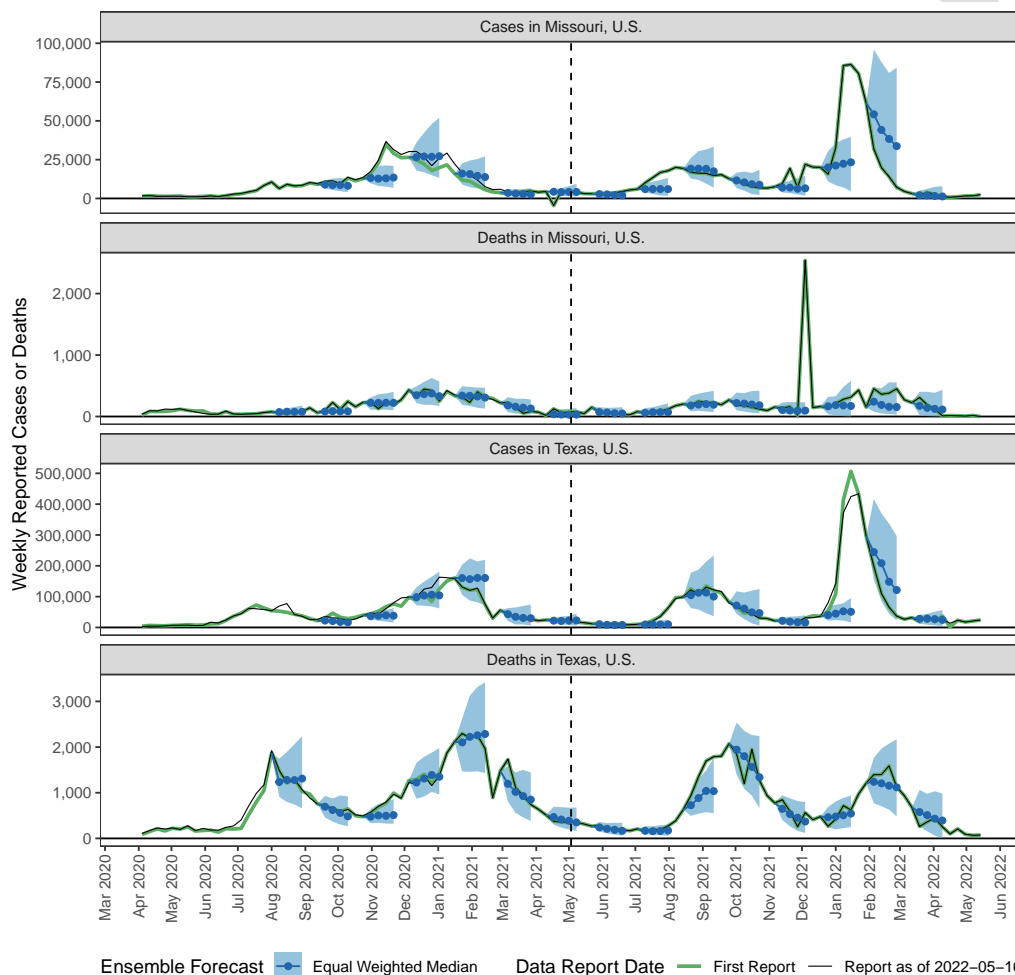


Figure 2: Weekly reported cases and deaths and example equally weighted median ensemble forecasts (predictive median and 95% interval) for selected U.S. states. Forecasts were produced each week, but for legibility, only forecasts originating from every sixth week are displayed. Data providers occasionally change initial reports (green lines) leading to revised values (black lines). Vertical dashed lines indicate the start of the prospective ensemble evaluation phase.

1
2
3 submission was required. The Forecast Hubs enforced other validation criteria, including that
4 predictions of incident cases and deaths were non-negative and predictive quantiles were properly
5 ordered across quantile levels.
6
7

8 9 10 *2.3. Notation*

11 We denote the reported number of cases or deaths for location l and week t by $y_{l,t}$. A single
12 predictive quantile from component forecaster m is denoted by $q_{l,s,t,k}^m$, where s indexes the week
13 the forecast was created, t indexes the target week of the forecast, and k indexes the quantile
14 level. The forecast horizon is the difference between the target date t and the forecast date s .
15 There are a total of $K = 7$ quantile levels for forecasts of cases in the U.S., and $K = 23$ quantile
16 levels otherwise. The quantile levels are denoted by τ_k (e.g., if $\tau_k = 0.5$ then $q_{l,s,t,k}^m$ is a predictive
17 median). We collect the full set of predictive quantiles for a single model, location, forecast date,
18 and target date in the vector $q_{l,s,t,1:K}^m$. We denote the total number of available forecasters by M ;
19 this changes for different locations and weeks, but we suppress that in the notation.
20
21
22
23
24
25
26
27
28

29 *2.4. Baseline forecaster*

30 In the results below, many comparisons are made with reference to an epidemiologically naive
31 baseline forecaster that projects forward the most recent observed value with growing uncertainty
32 at larger horizons. This baseline forecaster was a random walk model on weekly counts of cases or
33 deaths, with $Y_{l,t} | Y_{l,t-1} = Y_{l,t-1} + \varepsilon_{l,t}$. The model used a non-parametric estimate of the distribution
34 of the innovations $\varepsilon_{l,t}$ based on the observed differences in weekly counts $d_{l,s} = y_{l,s} - y_{l,s-1}$ over
35 all past weeks s for the specified location l . Predictive quantiles were based on the quantiles of
36 the collection of these differences and their negations, using the default method for calculating
37 quantiles in R. The inclusion of negative differences ensured that the predictive distributions were
38 symmetric and the predictive median was equal to the most recent observed value. Forecasts
39 at horizons greater than one were obtained by iterating one-step-ahead forecasts. Any resulting
40 predictive quantiles that were less than zero were truncated to equal zero.
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

2.5. Evaluation metrics

To evaluate forecasts, we adopted the *weighted interval score* (WIS) (Bracher et al., 2021a). Let $q_{1:K}$ be predictive quantiles for the observed quantity y . The WIS is calculated as

$$\text{WIS}(q_{1:K}, y) = \frac{1}{K} \sum_{k=1}^K 2 \{ \mathbb{1}_{(-\infty, q_k]}(y) - \tau_k \} (q_k - y),$$

where $\mathbb{1}_{(-\infty, q_k]}(y)$ is the indicator function that takes the value 1 when $y \in (-\infty, q_k]$ and 0 otherwise. This is a negatively-oriented proper score, meaning that negative scores are better and its expected value according to a given data generating process is minimized by reporting the predictive quantiles from that process. WIS was designed as a discrete approximation to the continuous ranked probability score, and is equivalent to pinball loss, which is commonly used in quantile regression (Bracher et al., 2021a). We note that some other commonly used scores such as the logarithmic score and the continuous ranked probability score are not suitable for use with predictive distributions that are specified in terms of a set of predictive quantiles, since a full predictive density or distribution function is not directly available (see Supplemental Section 3 for further discussion).

To compare the skill of forecasters that submitted different subsets of forecasts, we used *relative WIS*, as done in Cramer et al. (2022). The ensemble forecasters developed and evaluated in this manuscript provided all relevant forecasts; missingness pertains only to the component forecasters, and in the present work the relative WIS is primarily used to summarize component forecaster skill as an input to some of the trained ensemble methods described below. Let \mathcal{I} denote a set of combinations of location l and forecast creation date s over which we desire to summarize model performance, and $\mathcal{I}_{m,m'} \subseteq \mathcal{I}$ be the subset of those locations and dates for which both models m and m' provided forecasts through a forecast horizon of at least four weeks. The relative WIS of model m over the set \mathcal{I} is calculated as

$$\text{rWIS}_{\mathcal{I}}^m = \frac{\theta^m}{\theta^{\text{baseline}}}, \text{ where}$$

$$\theta^m = \left(\frac{\prod_{m'=1}^M (4 \cdot |\mathcal{I}_{m,m'}|)^{-1} \sum_{(l,s) \in \mathcal{I}_{m,m'}} \sum_{t=s+1}^{s+4} \text{WIS}(q_{l,s,t,1:K}^m, y_{l,t})}{(4 \cdot |\mathcal{I}_{m,m'}|)^{-1} \sum_{(l,s) \in \mathcal{I}_{m,m'}} \sum_{t=s+1}^{s+4} \text{WIS}(q_{l,s,t,1:K}^{m'}, y_{l,t})} \right)^{\frac{1}{M}}.$$

In words, we computed the ratio of the mean WIS scores for model m and each other model m' , averaging across the subset of forecasts shared by both models. θ^m was calculated as the geometric

mean of these pairwise ratios of matched mean scores, and summarized how model m did relative to all other models on the forecasts they had in common. These geometric means were then scaled such that the baseline forecaster had a relative WIS of 1; a relative WIS less than 1 indicated forecast skill that was better than the baseline model. We note that if no forecasts were missing, $\mathcal{I}_{m,m'}$ would be the same for all model pairs, so that the denominators of each θ^m and of θ^{baseline} would cancel when normalizing relative to the baseline and the relative WIS for model m would reduce to the mean WIS for model m divided by the mean WIS for the baseline model. We used the geometric mean to aggregate across model pairs to match the convention set in Cramer et al. (2022), but this detail is not critical: Supplemental Figure 5 illustrates that the relative WIS changes very little if an arithmetic mean is used instead.

We also assessed probabilistic calibration of the models with the one-sided coverage rates of predictive quantiles, calculated as the proportion of observed values that were less than or equal to the predicted quantile value. For a well-calibrated model, the empirical one-sided coverage rate is equal to the nominal quantile level. A method that generates conservative two-sided intervals would have an empirical coverage rate that is less than the nominal rate for quantile levels less than 0.5 and empirical coverage greater than the nominal rate for quantile levels greater than 0.5.

2.6. Ensemble model formulations

All of the ensemble formulations that we considered obtain a predictive quantile at level k by combining the component forecaster predictions at that quantile level:

$$q_{l,s,t,k}^{\text{ens}} = f(q_{l,s,t,k}^1, \dots, q_{l,s,t,k}^M).$$

We conceptually organize the ensemble methods considered according to two factors. First, *trained* ensemble methods use the past performance of the component forecasters to select a subset of components for inclusion in the ensemble and/or assign the components different weights, whereas *untrained* methods assign all component forecasters equal weight. Second, we varied the robustness of the combination function f to outlying component forecasts. Specifically, we considered methods based on either a (weighted) mean, which can be sensitive to outlying forecasts, or a (weighted) median, which may be more robust to these outliers. The weighted mean calculates the ensemble quantiles as

$$q_{l,s,t,k}^{\text{ens}} = \sum_{m=1}^M w_s^m q_{l,s,t,k}^m.$$

The weighted median is defined to be the smallest value q for which the combined weight of all component forecasters with predictions less than or equal to q is at least 0.5; the ensemble forecast quantiles are calculated as:

$$q_{l,s,t,k}^{\text{ens}} = \inf \left\{ q \in \mathbb{R} : \sum_{m=1}^M w_s^m \mathbb{1}_{(-\infty, q]}(q_{l,s,t,k}^m) \geq 0.5 \right\}.$$

In practice, we used the implementation of the weighted median in the `matrixStats` package for R, which linearly interpolates between the central weighted sample quantiles (Bengtsson, 2020). Graphically, these ensembles can be interpreted as computing a horizontal mean or median of the CDFs of component forecasters (Supplemental Figure 7).

In trained ensemble methods that weight the component forecasters, the weights were calculated as a sigmoidal transformation of the forecasters' relative WIS (see Section 2.5) over a rolling window of weeks leading up to the ensemble forecast date s , denoted by rWIS_s^m :

$$w_s^m = \frac{\exp(-\theta_s \cdot \text{rWIS}_s^m)}{\sum_{m'=1}^M \exp(-\theta_s \cdot \text{rWIS}_s^{m'})}.$$

This formulation requires estimating the nonnegative parameter θ_s , which was updated each week. If $\theta_s = 0$, the procedure reduces to an equal weighting scheme. However, if θ_s is large, better-performing component forecasters (with low relative WIS scores) are assigned higher weight. We selected θ_s by using a grid search to optimize the weighted interval score of the ensemble forecast over the training window, summing across all locations and relevant target weeks on or before time s :

$$\theta_s = \arg \min_{\theta} \sum_l \sum_{r=s-1}^{s-a} \sum_{t=r+1}^{\min(r+4,s)} \text{WIS}(q_{l,r,t,1:K}^{\text{ens},\theta}, y_{l,t}).$$

The size of the training window, a , is a tuning parameter that must be selected; we considered several possible values during model development, discussed further below. In a post hoc analysis, we considered regularizing the weights by setting a limit on the weight that could be assigned to any one model. We implemented this regularization strategy by restricting the grid of values for θ_s to those values for which the largest component forecaster weight was less than the maximum weight limit.

In this parameterization, the component forecaster weights are by construction nonnegative and sum to 1. When forecasts were missing for one or more component forecasters in a particular location and forecast date, we set the weights for those forecasters to 0 and renormalized the weights for the remaining forecasters so that they summed to 1.

1
2
3
4 Some trained ensembles that we considered used a preliminary component selection step, where
5 the top few individual forecasters were selected for inclusion in the ensemble based on their relative
6 WIS during the training window. The number of component forecasters selected is a tuning
7 parameter that we explored during model development. This component selection step may be
8 used either in combination with the continuous weighting scheme described above, or with an
9 equally-weighted combination of selected forecasters. Throughout the text below, we use the term
10 “trained” ensemble to refer generically to a method that uses component selection and/or weighting
11 based on historical component forecaster performance.
12
13

14
15
16
17 There are many other weighted ensembling schemes that could be formulated. For example,
18 separate weights could be estimated for different forecast horizons, for different quantile levels, or for
19 subsets of locations. As another example, the weights could be estimated by directly minimizing
20 the WIS associated with look-ahead ensemble forecasts (Taylor & Taylor, 2021). We explored
21 these and other ideas during model development, but our analyses did not show them to lead to
22 substantial gains, and thus we settled on the simpler weighting schemes presented above. Further
23 discussion of alternative schemes is deferred to the supplement.
24
25
26
27
28
29
30

31 *2.7. Data and code accessibility*

32

33 All component model forecasts and code used for fitting ensemble models and conducting the
34 analyses presented in this manuscript are available in public GitHub repositories (Cramer et al.,
35 2021b; Ray, 2020, 2021).
36
37
38
39

40 **3. Results**

41

42 We discuss the decisions that we made during model development in Section 3.1 before turning
43 to a more focused discussion of the impact on ensemble forecast skill of using robust or non-
44 robust combination mechanisms in Section 3.2, and trained or untrained methods in Section 3.3.
45 Section 3.4 presents a post hoc evaluation of a variation on ensemble methods that regularizes the
46 component forecaster weights. Results for the evaluation using forecasts in Europe are presented
47 in Section 3.5.
48
49
50
51
52

53 Throughout this section, scores were calculated using the ground truth data that were available
54 as of May 16, 2022 unless otherwise noted. This allowed five weeks of revisions to accrue between
55 the last target end date that was evaluated and the date of the data used for evaluation. When
56
57
58
59

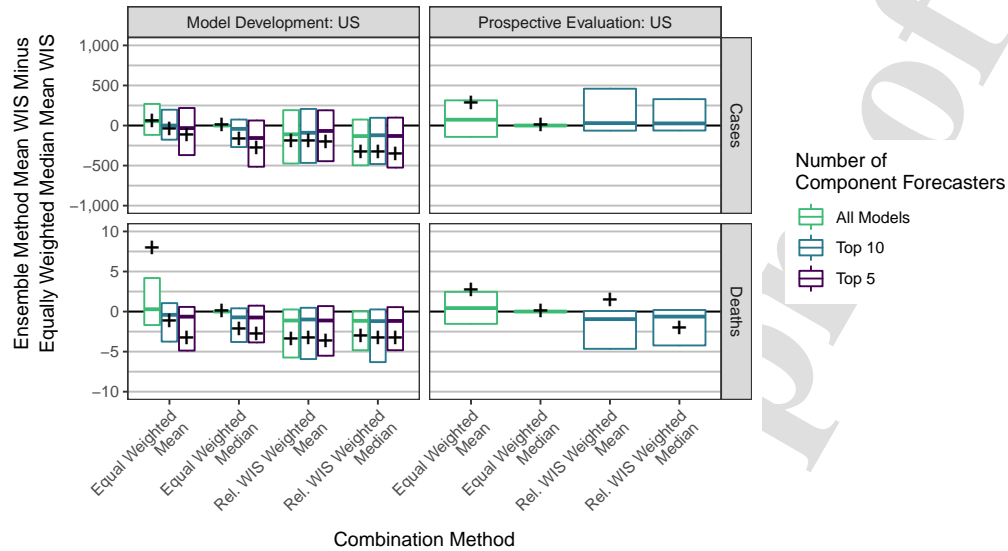
1
2
3 reporting measures of forecast skill, we dropped forecasts for which the corresponding reported
4 value of weekly cases or deaths was negative. This could occur when an error in data collection
5 was identified and corrected, or when the definition of a reportable case or death was changed.
6
7 We included scores for all other outlying and revised data in the primary analysis because it was
8 difficult to define objective standards for what should be omitted. However, a supplemental analysis
9 indicated that the results about the relative performance of different ensemble methods were not
10 sensitive to these reporting anomalies (Supplemental Section 5.4, Supplemental Figures 14 through
11 16).

12 13 14 15 16 17 18 19 *3.1. Model development*

20
21 During model development, we evaluated many variations on trained ensemble methods. In
22 these comparisons we take the equally weighted median ensemble as a reference approach because
23 this is the method used for the production ensemble produced by the U.S. Hub during most of
24 the time that we were running these experiments. As measured by mean WIS over the model
25 development phase, the equally weighted median ensemble was better than the equally weighted
26 mean ensemble, but both were outperformed by the trained ensemble variations using component
27 forecaster selection and/or weighting (Figure 3). The weighted approaches had stable performance
28 no matter how many component forecasters were included. Approaches using an equally weighted
29 combination of selected component forecasters were generally better only when top-performing
30 component forecasters were included.

31
32 We also considered varying other tuning parameters such as the length of the training window
33 and whether component forecaster weights were shared across different quantile levels or across
34 forecast horizons. However, we did not find strong and consistent gains in performance when
35 varying these other factors (Supplemental Figures 17 through 22). Finally, we evaluated other
36 possible formulations of weighted ensembles, with weights that were not directly dependent on the
37 relative WIS of the component forecasters but were instead estimated by optimizing the look-ahead
38 ensemble WIS over the training set. As measured by mean WIS, the best versions of these other
39 variations on weighted ensembles had similar performance to the best versions of the relative WIS
40 weighted median considered in the primary analysis. However, they were more sensitive to settings
41 like the number of component forecasters included and the training set window size (Supplemental
42 Figures 17 and 18).

(a) Weighted Interval Scores



(b) One-sided Quantile Coverage Rates

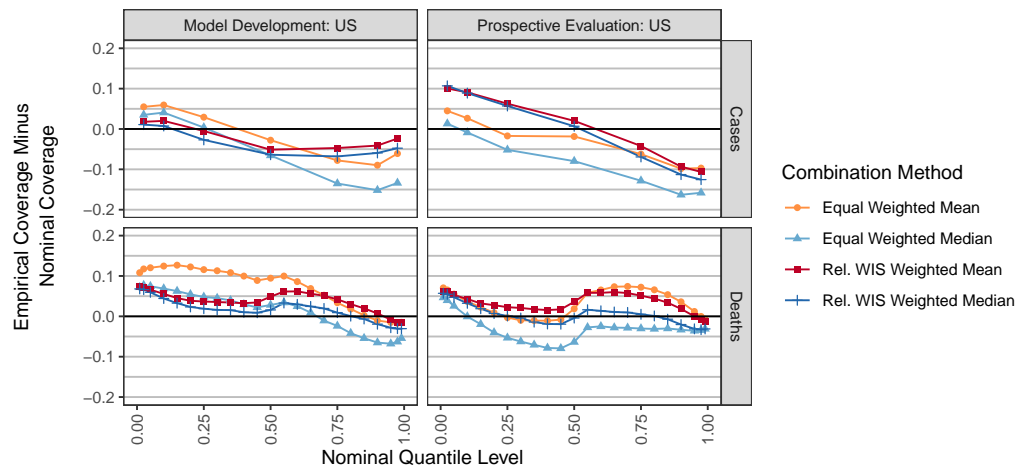


Figure 3: Performance measures for ensemble forecasts of weekly cases and deaths at the state level in the U.S. In panel (a) the vertical axis is the difference in mean WIS for the given ensemble method and the equally weighted median ensemble. Boxes show the 25th percentile, 50th percentile, and 75th percentile of these differences, averaging across all locations for each combination of forecast date and horizon. For legibility, outliers are suppressed here; Supplemental Figure 8 shows the full distribution. A cross is displayed at the difference in overall mean scores for the specified combination method and the equally weighted median averaging across all locations, forecast dates, and horizons. Large mean score differences of approximately 2,005 and 2,387 are suppressed for the Rel. WIS Weighted Mean and Rel. WIS Weighted Median ensembles respectively in the prospective phase forecasts of cases. A negative value indicates that the given method outperformed the equally weighted median. The vertical axis of panel (b) shows the probabilistic calibration of the ensemble forecasts through the one-sided empirical coverage rates of the predictive quantiles. A well-calibrated forecaster has a difference of 0 between the empirical and nominal coverage rates, while a forecaster with conservative (wide) two-sided intervals has negative differences for nominal quantile levels less than 0.5 and positive differences for quantile levels greater than 0.5.

1
2
3
4 Based on these results, on May 3, 2021 we selected the relative WIS weighted ensemble vari-
5 ations for use in the prospective evaluation, as these methods had similar mean WIS as the best
6 of the other variations considered, but were more consistent across different training set window
7 sizes and numbers of component forecasters included. We used intermediate values for these tun-
8 ing parameter settings, including 10 component forecasters with a training set window size of 12
9 weeks. We also included the equally weighted mean and median of all models in the prospective
10 evaluation as reference methods. The following sections give a more detailed evaluation of these
11 selected methods, describing how they performed during both the model development phase and
12 the prospective evaluation phase.
13
14
15
16
17
18

19 *3.2. Comparing robust and non-robust ensemble methods*

20
21
22 We found that for equally weighted ensemble approaches, robust combination methods were
23 helpful for limiting the effects of outlying component forecasts. For most combinations of evaluation
24 phase (model development or prospective evaluation) and target variable (cases or deaths), the
25 equally weighted median had better mean and worst-case WIS than the equally weighted mean,
26 often by a large margin (Figure 3, Supplemental Figure 8). Results broken down by forecast date
27 show that the methods achieved similar scores most of the time, but the equally weighted mean
28 ensemble occasionally had serious failures (Supplemental Figure 10). These failures were generally
29 associated with instances where a component forecaster issued extreme, outlying forecasts, e.g.,
30 forecasts of deaths issued the week of February 15th in Ohio (Figure 1).
31
32
33
34
35
36
37

38
39 There were fewer consistent differences between the trained mean and trained median ensemble
40 approaches. This suggests that both trained approaches that we considered had similar robustness
41 to outlying forecasts (if the outliers were produced by component forecasters that were down
42 weighted or not selected for inclusion due to poor historical performance) or sensitivity to outlying
43 forecasts (if they were produced by component forecasters that were selected and given high weight).
44
45
46
47

48
49 Panel (b) of Figure 3 summarizes probabilistic calibration of the ensemble forecasts with one-
50 sided quantile coverage rates. The median-based ensemble approaches generally had lower one-
51 sided quantile coverage rates than the mean-based approaches, indicating a downward shift of
52 the forecast distributions. This was associated with poorer probabilistic calibration for forecasts of
53 cases, where the ensemble forecast distributions tended to be too low. For forecasts of deaths, which
54 were better centered but tended to be too narrow, the calibration of the median-based methods was
55
56
57
58
59

1
2
3 not consistently better or worse than the calibration of the corresponding mean-based methods.
4
5

6 *3.3. Comparing trained and untrained ensemble methods*

7

8 Averaging across all forecasts for incident cases and deaths in the model development phase, the
9 weighted median was better than the equally weighted median and the weighted mean was better
10 than the equally weighted mean (Figure 3). However, in the prospective evaluation, the trained
11 methods showed improved mean WIS relative to untrained methods when forecasting deaths, but
12 were worse when forecasting cases. In general, the trained ensembles also came closer to matching
13 the performance of a post hoc weighted mean ensemble for deaths than for cases (Figures 4 and
14 5). This post hoc weighted mean ensemble estimated the optimal weights for each week after the
15 forecasted data were observed; it would not be possible to use this method in real time, but it gives
16 a bound on the ensemble forecast skill that can be achieved using quantile averaging.
17
18

19 We believe that this difference in the relative performance of trained and untrained ensemble
20 methods for cases and deaths is primarily due to differences in component model behavior for
21 forecasting cases and deaths. A fundamental difference between these outcomes is that cases are
22 a leading indicator relative to deaths, so that trends in cases in the recent past may be a helpful
23 input for forecasting deaths—but there are not clear candidates for a similar leading indicator for
24 cases (e.g., see McDonald et al. (2021) for an investigation of some possibilities that were found
25 to yield only modest and inconsistent improvements in forecast skill). Indeed, the best models for
26 forecasting mortality generally do use previously reported cases as an input to forecasting (Cramer
27 et al., 2022), and it has previously been noted that deaths are an easier target to forecast than
28 cases (Reich et al., 2021 [Online]; Bracher et al., 2021b). This is reflected in the performance of
29 trained ensembles, which were often able to identify a future change in direction of trends when
30 forecasting deaths, but generally tended to predict a continuation of recent trends when forecasting
31 cases (Supplemental Section 7, Supplemental Figures 25 and 26). An interpretation of this is that
32 the component forecasters with the best record of performance for forecasting deaths during the
33 training window were able to capture changes in trend, but the best component forecasters for
34 forecasting cases were often simply extrapolating recent trends. While all ensemble methods tended
35 to “overshoot” at local peaks in weekly incidence, this tendency was more pronounced for forecasts
36 of cases than for forecasts of deaths—and training tended to exacerbate the tendency to overshoot
37 when forecasting cases, but to mitigate this tendency when forecasting deaths (Supplemental Figure
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

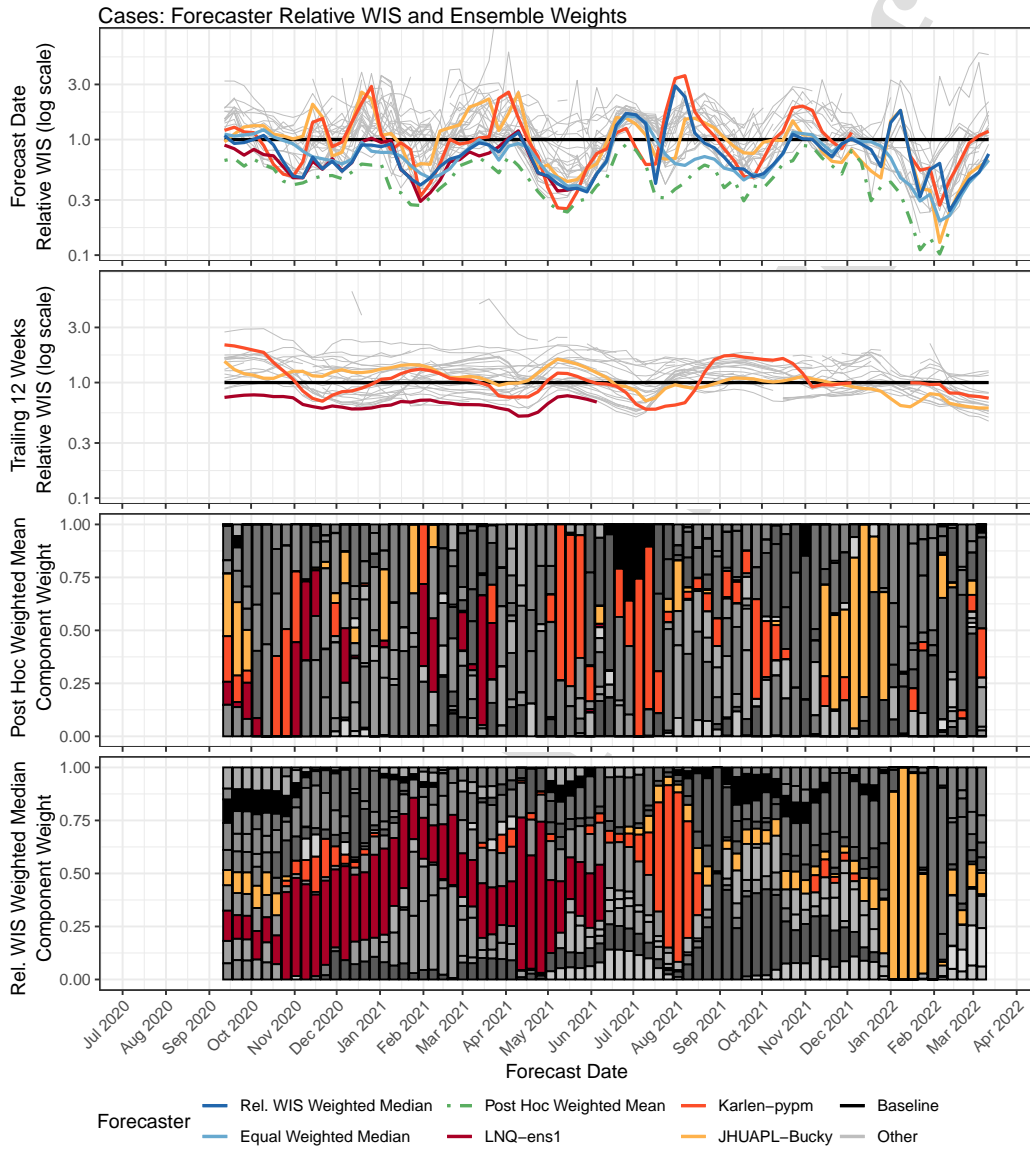


Figure 4: Performance of weekly case forecasts from component forecasters and selected ensembles, along with component forecaster weights. Component forecasters that were given high weight at key times are highlighted. The top row shows the relative WIS of forecasts made each week. The second row shows the relative WIS over the 12 weeks before the forecast date, for forecasts of quantities that were observed by the forecast date. These scores, which are used to compute the component weights in the relative WIS weighted median ensemble, are calculated using data available as of the forecast date. The third row shows component forecaster weights for the post hoc weighted mean ensemble, and the bottom row shows the component model weights for the relative WIS weighted median ensemble; each component forecaster is represented with a different color. Over the time frame considered, 31 distinct component forecasters were included in this top-10 ensemble.

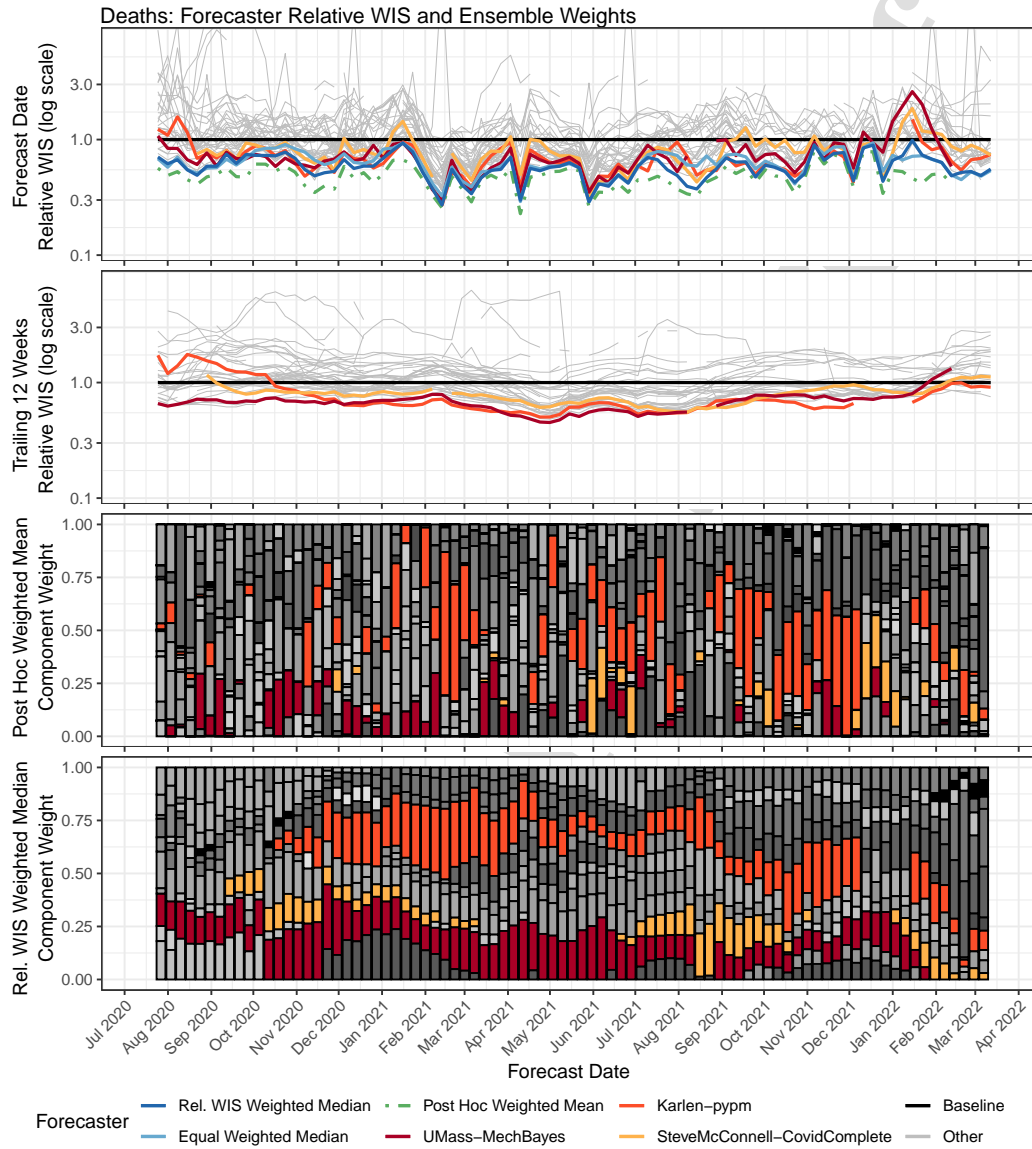


Figure 5: Performance of weekly death forecasts from component forecasters and selected ensembles, along with component forecaster weights. Component forecasters that were given high weight at key times are highlighted. The top row shows the relative WIS of forecasts made each week. The second row shows the relative WIS over the 12 weeks before the forecast date, for forecasts of quantities that were observed by the forecast date. These scores, which are used to compute the component weights in the relative WIS weighted median ensemble, are calculated using data available as of the forecast date. The third row shows component forecaster weights for the post hoc weighted mean ensemble, and the bottom row shows the component model weights for the relative WIS weighted median ensemble; each component forecaster is represented with a different color. Over the time frame considered, 34 distinct component forecasters were included in this top-10 ensemble.

1
2
3
4 25).

5 Another difference in component behavior when forecasting cases and deaths is illustrated in
6 Figures 4 and 5, which explore the relationships between component forecaster performance and
7 the relative performance of trained and untrained ensemble methods in more detail. For deaths,
8 the trained ensemble was able to identify and upweight a few component forecasters that had
9 consistently good performance (e.g., Karlen-pypm and UMass-MechBayes). This led to consistently
10 strong performance of the trained ensemble; it was always among the best models contributing to
11 the U.S. Hub, and was better than the equally weighted median ensemble in nearly every week.
12

13
14
15
16
17 For cases, the trained ensemble also had strong performance for many months when the LNQ-
18 ens1 forecaster was contributing to the U.S. Hub. However, when LNQ-ens1 stopped contributing
19 forecasts in June 2021, the trained ensemble shifted to weighting Karlen-pypm, which had less
20 stable performance for forecasting cases. During July 2021, Karlen-pypm was the only forecaster
21 in the U.S. Hub that predicted rapid growth at the start of the Delta wave, and it achieved the best
22 relative WIS by a substantial margin at that time. However, that forecaster predicted continued
23 growth as the Delta wave started to wane and it had the worst relative WIS a few weeks later. A
24 similar situation occurred during the Omicron wave in January 2022, when the JHUAPL-Bucky
25 model was one of a small number of forecasters that captured the rise at the beginning of the wave,
26 but it then overshot near the peak. In both of these instances, the post hoc weighting would have
27 assigned a large amount of weight to the forecaster in question at the start of the wave, when it was
28 uniquely successful at identifying rising trends in cases—but then shifted away from that forecaster
29 as the peak neared. Trained ensembles that estimated weights based on past performance suffered,
30 as they started to upweight those component forecasters just as their performance dropped. This
31 recurring pattern highlights the challenge that nonstationary component forecaster performance
32 presents for trained ensembles. Reinforcing this point, we note that in the post hoc weighted
33 mean ensemble, the component forecaster weights are only weakly autocorrelated (Figures 4 and
34 5, Supplemental Figure 27), again suggesting that an optimal weighting may require frequently
35 changing component weights to adapt to nonstationary performance.
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50

51 During the model development phase, the trained ensembles had better probabilistic calibration
52 than their equally weighted counterparts (Figure 3 panel (b)). During prospective evaluation, the
53 trained median ensemble had generally higher one-sided coverage rates, corresponding to better
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3 calibration in the upper tail but slightly worse calibration in the lower tail. The trained mean
4 ensemble had slightly better calibration than the equally weighted mean when forecasting deaths
5 in the prospective evaluation phase, but inconsistent gains across different quantile levels when
6 forecasting cases. Supplemental Figures 12 and 13 show that the widths of 95% prediction intervals
7 from both the equally weighted median ensemble and the relative WIS weighted median ensemble
8 tended to rank near the middle of the widths of 95% prediction intervals from the component
9 forecasters. This can be interpreted as an advantage if we are concerned about the possible influence
10 of component forecasters with very narrow or very wide prediction intervals. However, it can also
11 be viewed as a disadvantage, particularly if improved calibration could have been realized if the
12 prediction intervals were wider. We return to this point in the discussion.
13
14
15
16
17
18
19
20
21

22 *3.4. Post hoc evaluation of weight regularization*

23

24 Motivated by the assignment of large weights to some component forecasters in the trained
25 ensembles for cases (Figure 4), in January 2022 we conducted a post hoc evaluation of trained
26 ensembles that were regularized by imposing a limit on the weight that could be assigned to
27 any one component forecaster (see Section 2.6). In this evaluation, we constructed relative WIS
28 weighted median ensemble forecasts for all historical forecast dates up through the week of January
29 3, 2022. These ensemble fits included the top 10 component forecasters and were trained on a
30 rolling window of the 12 most recent forecast dates, matching the settings that were selected for
31 the prospective analysis. We considered six values for the maximum weight limit: 0.1, 0.2, 0.3,
32 0.4, 0.5, and 1.0. A weight limit of 1.0 corresponds to the unregularized method considered in the
33 prospective evaluation, and a weight limit of 0.1 corresponds to an equally weighted median of the
34 top ten forecasters, which was previously considered during the model development phase.
35
36
37
38
39
40
41
42
43

44 For both cases and deaths, the results of this analysis indicate that a weight limit as low as
45 0.1 was unhelpful (Figure 6). When forecasting deaths, this regularization strategy had limited
46 impact on the trained ensemble performance as long as the maximum weight limit was about 0.3 or
47 higher, which is consistent with the fact that the trained ensembles for deaths rarely assigned a large
48 weight to one model (Figure 5). However, when forecasting cases, the regularization resulted in
49 large improvements in mean WIS, with the best WIS at limits near 0.2 or 0.3. These improvements
50 were concentrated in short periods near local peaks in the epidemic waves (Supplemental Figure
51 28). For both cases and deaths, smaller limits on the maximum weight were associated with a
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3 slight reduction in the empirical coverage rates of 95% prediction intervals. Based on these results,
4 the U.S. Hub used a weight limit of 0.3 in trained ensemble forecasts starting in January 2022.
5
6
7

8 *3.5. Results in the European application*

9

10 Figure 7 summarizes weighted interval scores and calibration for the four selected ensemble
11 methods when applied prospectively to forecast data collected in the European Forecast Hub.
12 Consistent with what we observed for the U.S. above, the equally weighted median ensemble was
13 generally better than the equally weighted mean. However, in the European evaluation, the trained
14 methods had worse performance than the equally weighted median for forecasting both cases and
15 deaths.
16
17
18
19

20 In a post hoc exploratory analysis, we noted that patterns of missingness in forecast submissions
21 are quite different in the U.S. and in Europe (Figure 8, Supplemental Figures 29 through 36). In the
22 U.S. Hub, nearly all models submit forecasts for all of the 50 states, and many additionally submit
23 forecasts for at least one of the District of Columbia and territories. However, in the European Hub,
24 roughly half of contributing models submit forecasts for only a small number of locations. Because
25 the trained ensembles selected for prospective evaluation select the top 10 individual forecasters
26 by relative WIS, this means that in practice the trained ensembles only included a few component
27 forecasters for many locations in Europe.
28
29
30
31
32
33
34
35

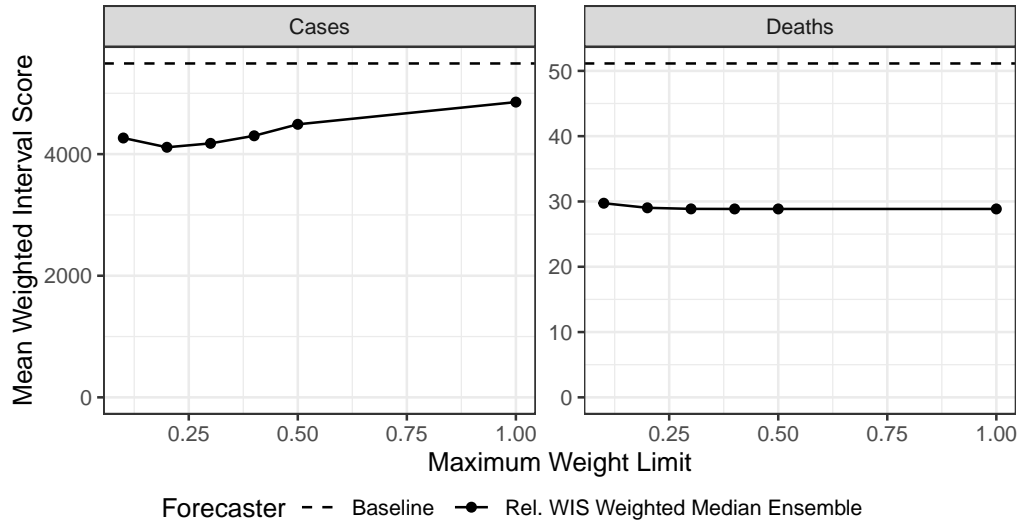
36 **4. Discussion**

37
38

39 In this work, we have documented the analyses that have informed the selection of methods
40 employed by the official U.S. Hub ensemble that is used by the CDC for communication with
41 public health decision makers and the public more generally. In this context, our preference is for
42 methods that have stable performance across different locations and different points in time, and
43 good performance on average.
44
45
46
47

48 Our most consistent finding is that robust ensemble methods (i.e., based on a median) are
49 helpful because they are more stable in the presence of outlying forecasts than methods using a
50 mean. Ensemble methods based on means have repeatedly produced extreme forecasts that are
51 dramatically misaligned with the observed data, but median-based approaches have not suffered
52 from this problem as much. This stability is of particular importance in the context of forecasts
53
54
55
56
57
58
59
60
61
62
63
64
65

(a) Ensemble WIS by maximum weight limit



(b) Ensemble 95% interval coverage rate by maximum weight limit

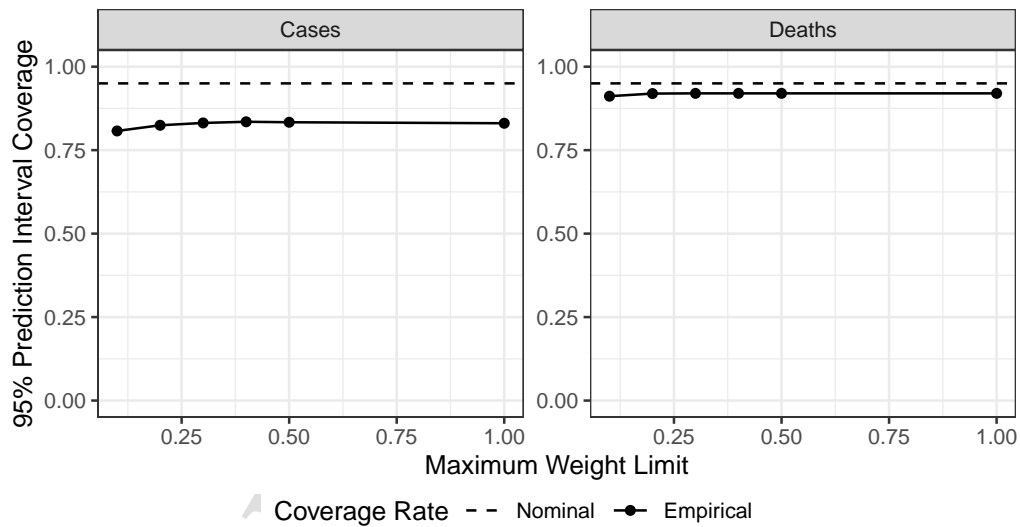
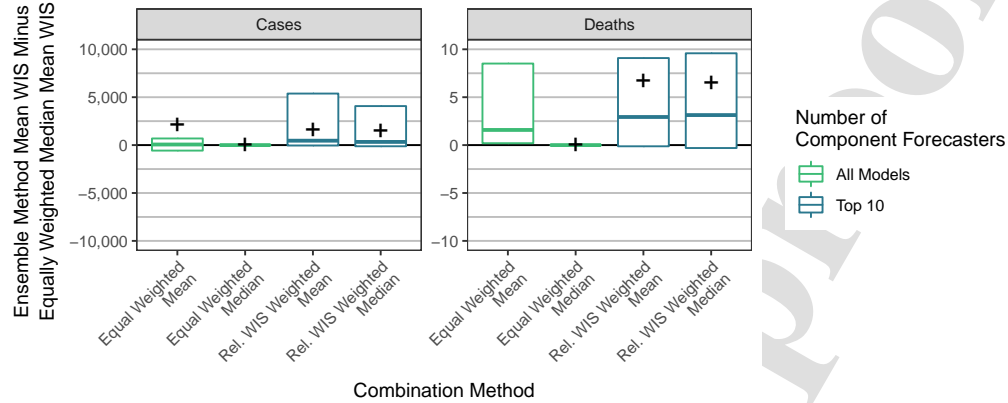


Figure 6: Mean WIS and 95% prediction interval coverage rates for relative WIS weighted median trained ensemble variations with varying sizes of a limit on the weight that could be assigned to any one model. In panel (a), the baseline forecaster is included as a reference. Results are for a post hoc analysis including forecast dates up to January 3, 2022.

(a) Weighted Interval Scores



(b) One-sided Quantile Coverage Rates

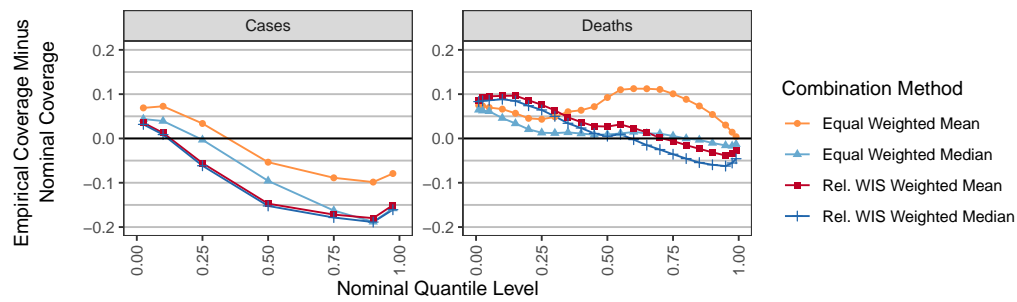
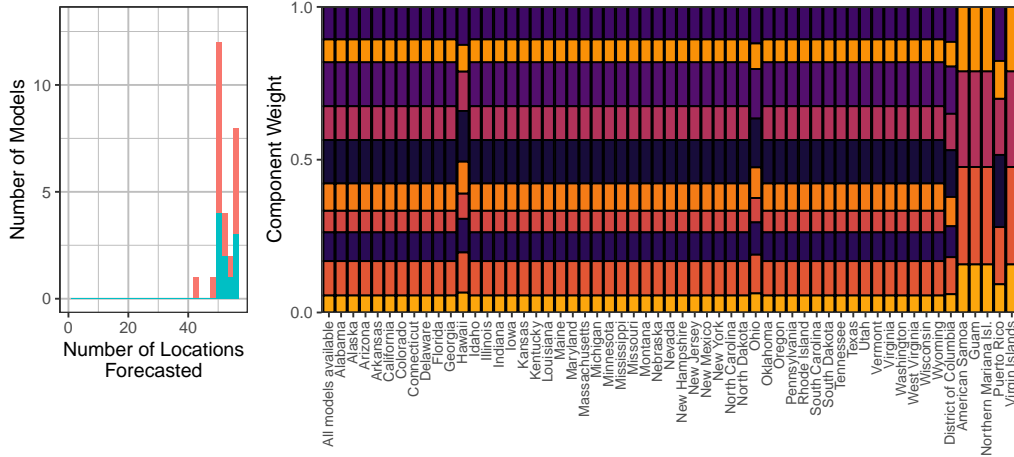


Figure 7: Performance measures for ensemble forecasts of weekly cases and deaths in Europe. In panel (a) the vertical axis is the difference in mean WIS for the given ensemble method and the equally weighted median ensemble. Boxes show the 25th percentile, 50th percentile, and 75th percentile of these differences, averaging across all locations for each combination of forecast date and horizon. For legibility, outliers are suppressed here; Supplemental Figure 9 shows the full distribution. A cross is displayed at the difference in overall mean scores for the specified combination method and the equally weighted median of all models, averaging across all locations, forecast dates, and horizons. A large mean score difference of approximately 666 is suppressed for the Equal Weighted Mean ensemble forecasts of deaths. A negative value indicates that the given method had better forecast skill than the equally weighted median. Panel (b) shows the probabilistic calibration of the forecasts through the one-sided empirical coverage rates of the predictive quantiles. A well-calibrated forecaster has a difference of 0 between the empirical and nominal coverage rates, while a forecaster with conservative (wide) two-sided intervals have negative differences for nominal quantile levels less than 0.5 and positive differences for quantile levels greater than 0.5.

(a) US: Number of locations forecasted per model and effective model weights per location



(b) EU: Number of locations forecasted per model and effective model weights per location

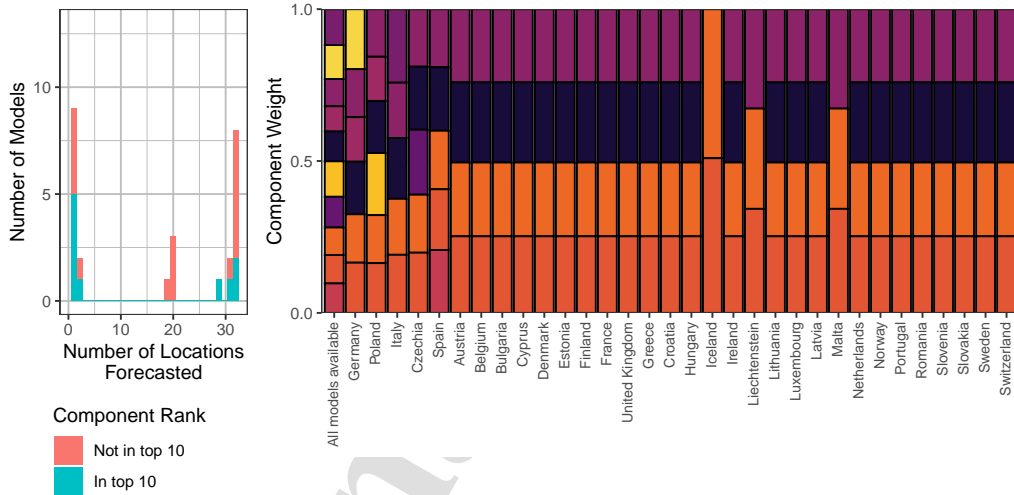


Figure 8: A comparison of the impacts of forecast missingness in the applications to the U.S. (panel (a)) and Europe (panel (b)). Within each panel, the histogram on the left shows the number of locations forecasted by each contributing forecaster in the week of October 11, 2021, colored by whether or not the forecaster was among the top 10 forecasters eligible for inclusion in the relative WIS weighted ensemble selected for prospective evaluation. The plot on the right shows the estimated weights that would be used if all of the top 10 models (each represented by a different color) were available for a given location (at left side), and the effective weights used in each location after setting the weights for models that did not provide location-specific forecasts to 0 and rescaling the other weights proportionally to sum to 1.

1
2
3 that will be used by public health decision makers. These observations informed our decision to
4 use an equally weighted median ensemble for the official U.S. Hub ensemble early on.
5
6

7 We have seen more mixed success for trained ensemble methods. Overall, trained ensemble
8 methods did well when they were able to identify and upweight component forecasters with stable
9 good performance, but struggled when component forecaster skill varied over time. In the U.S.,
10 trained ensembles have a long record of good performance when forecasting deaths, and the U.S.
11 Hub adopted the relative WIS weighted median ensemble as its official method for forecasting
12 deaths in November 2021. However, trained methods have been less successful at forecasting cases
13 in the U.S., both near peaks in weekly incidence (when they tend to overshoot) and at points
14 where the performance of the component forecasters is inconsistent. Additionally, the trained
15 methods we adopted did not translate well to a setting with a large number of missing component
16 forecasts, as in the European Hub. To preserve the prospective nature of our analyses, we have
17 not examined additional ensemble variations in the European application, but we hypothesize that
18 these problems might be mitigated by including all component forecasters rather than the top 10,
19 or by performing weight estimation separately in clusters of locations where the same component
20 forecasters are contributing. Allowing for different weights in different locations may also be an
21 effective strategy for addressing the impacts of differences in data availability and quality across
22 different locations.
23
24
25
26
27
28
29
30
31
32
33
34

35 In this manuscript, we have focused on relatively simple approaches to building ensemble fore-
36 casts. There are several opportunities for other directions that we have not considered here, and
37 the gap in performance between the ensemble methods we have considered and an ensemble us-
38 ing post hoc optimal weights indicates that there may still be room for improvement in ensemble
39 methods. In our view, the most central challenge for trained ensembles is the inconsistency of the
40 relative performance of many component forecasters, which may in turn be responsible for the lack
41 of strong short-term temporal correlation in the component forecaster weights that were estimated
42 by the post hoc weighted mean ensemble. For models with a relatively long history of performance
43 over multiple epidemic waves, we believe that the most promising approach to addressing this is
44 by using weights that depend on covariates like recent trends in incidence. This might allow the
45 ensemble to learn the conditions in which component forecasters have been more or less reliable,
46 and upweight models locally during phases similar to those in which they have done well in the
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3 past. Similar approaches have been used for other infectious disease systems such as influenza
4 in the past (e.g., Ray & Reich, 2018), but they used a substantial amount of training data over
5 multiple years.
6
7

8
9 There are several other possible directions for further exploration. We have addressed the chal-
10 lenge posed by outlying component forecasts by using median-based combination mechanisms, but
11 another approach would be to pre-screen the component forecasts and remove outlying forecasts.
12 This is a difficult task because there are times when weekly cases and deaths grow exponentially,
13 and occasionally only one or two models have captured this growth accurately (Supplemental
14 Figures 1 and 2). A component screening method would have to be careful to avoid screening
15 out methods that looked extreme relative to the data or other component forecasts, but in fact
16 accurately captured exponential growth (see Supplemental Section 1 for more discussion).
17
18

19
20 Another challenge is that the ensemble forecasts have not always been well calibrated. We are
21 actively developing approaches to address this by post hoc recalibration of the ensemble forecasts.
22 Another possible route forward would be to use a different method for ensemble construction. As we
23 discussed earlier, the ensemble methods that we have considered work by combining the predictions
24 from component forecasters at each quantile level, and therefore tend to have a dispersion that ranks
25 in the middle of the dispersions of the component forecasters. In contrast, an ensemble forecast
26 obtained as a distributional mixture of component forecasts would exhibit greater uncertainty
27 at times when the component forecasts disagreed with each other. However, such an approach
28 would be impacted by extreme component forecasts, and would likely require the development of
29 strategies for screening outlying forecasts as discussed above.
30
31

32
33 Additionally, our methods for constructing ensemble forecasts do not directly account for the
34 fact that some component forecasters are quite similar to each other and may provide redundant
35 information about the future of the pandemic. Ensembles generally benefit from combining diverse
36 component forecasters, and it could be helpful to encourage this—for example, by clustering
37 the forecasters and including a representative summary of the forecasts within each cluster as
38 the ensemble components. There are also related questions about the importance of different
39 component forecasters to ensemble skill; we plan to explore this direction in future work by using
40 tools such as the Shapley value to describe the contribution of individual components to the full
41 ensemble.
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 We have used the WIS and probabilistic calibration to measure the extent to which forecasts
5 are consistent with the eventually observed data. These summaries of performance are commonly
6 used and provide useful insights into forecast performance, but it is worth noting that they do
7 not necessarily reflect the utility of the forecasts for every particular decision-making context.
8
9 Aggregated summaries of performance, such as overall quantile coverage rates could obscure finer-
10 scale details—for instance, a method with good coverage rates on average could have high coverage
11 at times that are relatively unimportant and low coverage when it matters. Additionally, for some
12 public health decision-making purposes, one or another aspect of a forecast may be more important;
13 for example, some users may prioritize accurate assessments about when a new wave may begin,
14 but other users may find accurate forecasts of peak intensity to be more important. Our evaluation
15 metrics do not necessarily reflect the particular needs of those specific end users, and it is possible
16 that different ensemble methods would be more or less appropriate to generate forecasts that serve
17 different purposes.
18
19

20
21 Careful consideration and rigorous evaluation are required to support decisions about what
22 ensemble methods should be used for infectious disease forecasting. As we discussed earlier, to
23 obtain an accurate measure of a forecaster's performance, it is critical that the versions of ground
24 truth data that would have been available in real time are used for parameter estimation. This
25 applies as much to ensemble forecasters as it does to individual models. Additionally, it is important
26 to be clear about what methods development and evaluation were done retrospectively and what
27 forecasts were generated prospectively in real time. We believe that to avoid disruptions to public
28 health end users, a solid evidence base of stable performance in prospective forecasts should be
29 assembled to support a change in ensemble methods. We have followed these principles in this
30 work, and have followed the EPIFORGE guidelines in describing our analysis (Pollett et al. (2021);
31 Supplemental Section 11).
32
33

34
35 The COVID-19 pandemic has presented a unique challenge for infectious disease forecasting.
36 The U.S. and European Forecast Hubs have collected a wealth of forecasts from many contributing
37 teams—far more than have been collected in previous collaborative forecasting efforts for infectious
38 diseases such as influenza, dengue, and Ebola. These forecasts have been produced in real time to
39 respond to an emerging pathogen that has been one of the most serious public health crises in the
40 last century. This setting has introduced a myriad of modeling difficulties, from data anomalies due
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 to new reporting systems being brought online and changing case definitions, to uncertainty about
5 the fundamental epidemiological parameters of disease transmission, to rapidly changing social
6 factors such as implementation and uptake of non-pharmaceutical interventions. The behavior
7 of individual models in the face of these difficulties has in turn affected the methods that were
8 suitable for producing ensemble forecasts. We are hopeful that the lessons learned about infectious
9 disease forecasting will help to inform effective responses from the forecasting community in future
10 infectious disease crises.
11
12
13
14
15

16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65

- Agosto, A., Campmas, A., Giudici, P., & Renda, A. (2021). Monitoring COVID-19 contagion growth. *Statistics in Medicine*, *40*, 4150–4160. URL: <https://onlinelibrary.wiley.com/doi/10.1002/sim.9020>. doi:10.1002/sim.9020.
- Bartolucci, F., Pennoni, F., & Mira, A. (2021). A multivariate statistical approach to predict COVID-19 count data with epidemiological interpretation and uncertainty quantification. *Statistics in Medicine*, *40*, 5351–5372. URL: <https://onlinelibrary.wiley.com/doi/10.1002/sim.9129>. doi:10.1002/sim.9129.
- Bengtsson, H. (2020). *matrixStats: Functions that Apply to Rows and Columns of Matrices (and to Vectors)*. URL: <https://CRAN.R-project.org/package=matrixStats> R package version 0.57.0.
- Bracher, J., Ray, E. L., Gneiting, T., & Reich, N. G. (2021a). Evaluating epidemic forecasts in an interval format. *PLOS Computational Biology*, *17*, e1008618. doi:10.1371/journal.pcbi.1008618. arXiv:33577550.
- Bracher, J., Wolfram, D., Deuschel, J., Grgen, K., Ketterer, J. L., Ullrich, A., Abbott, S., Barbarossa, M. V., Bertsimas, D., Bhatia, S., Bodych, M., Bosse, N. I., Burgard, J. P., Castro, L., Fairchild, G., Fuhrmann, J., Funk, S., Gogolewski, K., Gu, Q., Heyder, S., Hotz, T., Kheifetz, Y., Kirsten, H., Krueger, T., Krymova, E., Li, M. L., Meinke, J. H., Michaud, I. J., Niedzielewski, K., Oaski, T., Rakowski, F., Scholz, M., Soni, S., Srivastava, A., Zieliski, J., Zou, D., Gneiting, T., & Schienle, M. (2021b). A pre-registered short-term forecasting study of COVID-19 in Germany and Poland during the second wave. *Nature Communications*, *12*, 5173. URL: <https://www.nature.com/articles/s41467-021-25207-0>. doi:10.1038/s41467-021-25207-0.
- Brooks, L. C., Farrow, D. C., Hyun, S., Tibshirani, R. J., & Rosenfeld, R. (2018). Nonmechanistic forecasts of seasonal influenza with iterative one-week-ahead distributions. *PLOS Computational Biology*, *14*, e1006134. URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6034894/>. doi:10.1371/journal.pcbi.1006134.
- Brooks, L. C., Ray, E. L., Bien, J., Bracher, J., Rumack, A., Tibshirani, R. J., & Reich, N. G. (2020 [Online]). Comparing ensemble approaches for short-term probabilistic COVID-19 forecasts in the U.S. International Institute of Forecasters blog. URL: <https://forecasters.org/blog/2020/10/28/comparing-ensemble-approaches-for-short-term-probabilistic-covid-19-forecasts-in-the-u-s/>.
- Busetti, F. (2017). Quantile aggregation of density forecasts. *Oxford Bulletin of Economics and Statistics*, *79*, 495–512.

- 1
2
3
4 CDC (2021). COVID-19 mathematical modeling. URL: <https://www.cdc.gov/coronavirus/2019-ncov/science/forecasting/mathematical-modeling.html>.
- 5
6
7 Claeskens, G., Magnus, J. R., Vasnev, A. L., & Wang, W. (2016). The forecast combination puzzle: A simple
8 theoretical explanation. *International Journal of Forecasting*, 32, 754–762. URL: <https://www.sciencedirect.com/science/article/pii/S0169207016000327>. doi:10.1016/j.ijforecast.2015.12.005.
- 9
10
11 Coln-Gonzalez, F. J., Bastos, L. S., Hofmann, B., Hopkin, A., Harpham, Q., Crocker, T., Amato, R., Ferrario, I.,
12 Moschini, F., James, S., Malde, S., Ainscoe, E., Nam, V. S., Tan, D. Q., Khoa, N. D., Harrison, M., Tsarouchi, G.,
13 Lumbroso, D., Brady, O. J., & Lowe, R. (2021). Probabilistic seasonal dengue forecasting in Vietnam: A modelling
14 study using superensembles. *PLOS Medicine*, 18, e1003542. URL: <https://journals.plos.org/plosmedicine/article?id=10.1371/journal.pmed.1003542>. doi:10.1371/journal.pmed.1003542. Publisher: Public Library of
15 Science.
- 16
17
18 Cramer, E. Y., Huang, Y., Wang, Y., Ray, E. L., Cornell, M., Bracher, J., Brennen, A., Rivadeneira, A. J. C.,
19 Gerding, A., House, K., Jayawardena, D., Kanji, A. H., Khandelwal, A., Le, K., Niemi, J., Stark, A., Shah, A.,
20 Wattanachit, N., Zorn, M. W., Reich, N. G., & U.S. COVID-19 Forecast Hub Consortium (2021a). The United
21 States COVID-19 Forecast Hub dataset. *medRxiv*, (p. 2021.11.04.21265886). URL: <https://www.medrxiv.org/content/10.1101/2021.11.04.21265886v1>. doi:10.1101/2021.11.04.21265886.
- 22
23
24 Cramer, E. Y., Huang, Y., Wang, Y., Ray, E. L., Cornell, M., Bracher, J., Brennen, A., Rivadeneira, A. J. C., Gerding,
25 A., House, K., Jayawardena, D., Kanji, A. H., Khandelwal, A., Le, K., Niemi, J., Stark, A., Shah, A., Wattanachit,
26 N., Zorn, M. W., Reich, N. G., & U.S. COVID-19 Forecast Hub Consortium (2021b). reichlab/covid19-forecast-
27 hub: release for zenodo, 20210816. URL: <https://zenodo.org/record/5208210>. doi:10.5281/zenodo.5208210.
- 28
29
30 Cramer, E. Y., Ray, E. L., Lopez, V. K., Bracher, J., Brennen, A., Rivadeneira, A. J. C., Gerding, A., Gneiting,
31 T., House, K. H., Huang, Y., Jayawardena, D., Kanji, A. H., Khandelwal, A., Le, K., Mhlemann, A., Niemi, J.,
32 Shah, A., Stark, A., Wang, Y., Wattanachit, N., Zorn, M. W., Gu, Y., Jain, S., Bannur, N., Deva, A., Kulkarni,
33 M., Merugu, S., Raval, A., Shingi, S., Tiwari, A., White, J., Abernethy, N. F., Woody, S., Dahan, M., Fox, S.,
34 Gaither, K., Lachmann, M., Meyers, L. A., Scott, J. G., Tec, M., Srivastava, A., George, G. E., Cegan, J. C.,
35 Dettwiller, I. D., England, W. P., Farthing, M. W., Hunter, R. H., Lafferty, B., Linkov, I., Mayo, M. L., Parno,
36 M. D., Rowland, M. A., Trump, B. D., Zhang-James, Y., Chen, S., Faraone, S. V., Hess, J., Morley, C. P.,
37 Salekin, A., Wang, D., Corsetti, S. M., Baer, T. M., Eisenberg, M. C., Falb, K., Huang, Y., Martin, E. T.,
38 McCauley, E., Myers, R. L., Schwarz, T., Sheldon, D., Gibson, G. C., Yu, R., Gao, L., Ma, Y., Wu, D., Yan,
39 X., Jin, X., Wang, Y.-X., Chen, Y., Guo, L., Zhao, Y., Gu, Q., Chen, J., Wang, L., Xu, P., Zhang, W., Zou,
40 D., Biegel, H., Lega, J., McConnell, S., Nagraj, V. P., Guertin, S. L., Hulme-Lowe, C., Turner, S. D., Shi, Y.,
41 Ban, X., Walraven, R., Hong, Q.-J., Kong, S. et al. (2022). Evaluation of individual and ensemble probabilistic
42 forecasts of COVID-19 mortality in the United States. *Proceedings of the National Academy of Sciences*, 119,
43 e2113561119. URL: <https://www.pnas.org/doi/abs/10.1073/pnas.2113561119>. doi:10.1073/pnas.2113561119.
44 arXiv:<https://www.pnas.org/doi/pdf/10.1073/pnas.2113561119>.
- 45
46
47 Dean, N. E., Pastore y Piontti, A., Madewell, Z. J., Cummings, D. A. T., Hitchings, M. D. T., Joshi, K., Kahn, R.,
48 Vespignani, A., Halloran, M. E., & Longini, I. M. (2020). Ensemble forecast modeling for the design of COVID-19
49 vaccine efficacy trials. *Vaccine*, 38, 7213–7216. URL: <https://www.sciencedirect.com/science/article/pii/S0264410X20300000>.
- 50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4 S0264410X20311919. doi:10.1016/j.vaccine.2020.09.031.
5 Dong, E., Du, H., & Gardner, L. (2020). An interactive web-based dashboard to track COVID-19 in real time.
6 *The Lancet Infectious Diseases*, 20, 533–534. URL: [https://www.thelancet.com/journals/laninf/article/](https://www.thelancet.com/journals/laninf/article/PIIS1473-3099(20)30120-1/fulltext)
7 [PIIS1473-3099\(20\)30120-1/fulltext](https://www.thelancet.com/journals/laninf/article/PIIS1473-3099(20)30120-1/fulltext). doi:10.1016/S1473-3099(20)30120-1. Publisher: Elsevier.
8 European Covid-19 Forecast Hub (2021). European Covid-19 Forecast Hub. URL: [https://covid19forecasthub.](https://covid19forecasthub.eu/)
9 [eu/](https://covid19forecasthub.eu/).
10
11 Gaba, A., Tsetlin, I., & Winkler, R. L. (2017). Combining interval forecasts. *Decision Analysis*, 14, 1–20.
12
13 Gneiting, T., Balabdaoui, F., & Raftery, A. E. (2007). Probabilistic forecasts, calibration and sharpness. *Journal of*
14 *the Royal Statistical Society: Series B (Statistical Methodology)*, 69, 243–268. URL: [https://onlinelibrary.](https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1467-9868.2007.00587.x)
15 [wiley.com/doi/abs/10.1111/j.1467-9868.2007.00587.x](https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1467-9868.2007.00587.x). doi:10.1111/j.1467-9868.2007.00587.x. eprint:
16 <https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1467-9868.2007.00587.x>.
17
18 Gneiting, T., & Raftery, A. E. (2005). Weather Forecasting with Ensemble Methods. *Science*, 310, 248–
19 249. URL: <https://www.science.org/doi/abs/10.1126/science.1115255>. doi:10.1126/science.1115255. Pub-
20 lisher: American Association for the Advancement of Science.
21
22 Gneiting, T., & Raftery, A. E. (2007). Strictly Proper Scoring Rules, Prediction, and Estimation. *Journal of the*
23 *American Statistical Association*, 102, 359–378. URL: <https://doi.org/10.1198/016214506000001437>. doi:10.
24 [1198/016214506000001437](https://doi.org/10.1198/016214506000001437). Publisher: Taylor & Francis eprint: <https://doi.org/10.1198/016214506000001437>.
25
26 Hora, S. C., Fransen, B. R., Hawkins, N., & Susel, I. (2013). Median aggregation of distribution functions. *Decision*
27 *Analysis*, 10, 279–291.
28
29 Johansson, M. A., Apfeldorf, K. M., Dobson, S., Devita, J., Buczak, A. L., Baugher, B., Moniz, L. J., Bagley,
30 T., Babin, S. M., Guven, E., Yamana, T. K., Shaman, J., Moschou, T., Lothian, N., Lane, A., Osborne, G.,
31 Jiang, G., Brooks, L. C., Farrow, D. C., Hyun, S., Tibshirani, R. J., Rosenfeld, R., Lessler, J., Reich, N. G.,
32 Cummings, D. A. T., Lauer, S. A., Moore, S. M., Clapham, H. E., Lowe, R., Bailey, T. C., Garca-Dez, M.,
33 Carvalho, M. S., Rod, X., Sardar, T., Paul, R., Ray, E. L., Sakrejda, K., Brown, A. C., Meng, X., Osoba, O.,
34 Vardavas, R., Manheim, D., Moore, M., Rao, D. M., Porco, T. C., Ackley, S., Liu, F., Worden, L., Convertino,
35 M., Liu, Y., Reddy, A., Ortiz, E., Rivero, J., Brito, H., Juarrero, A., Johnson, L. R., Gramacy, R. B., Cohen,
36 J. M., Mordecai, E. A., Murdock, C. C., Rohr, J. R., Ryan, S. J., Stewart-Ibarra, A. M., Weikel, D. P., Jutla, A.,
37 Khan, R., Poultney, M., Colwell, R. R., Rivera-Garca, B., Barker, C. M., Bell, J. E., Biggerstaff, M., Swerdlow,
38 D., Mier-y Teran-Romero, L., Forshey, B. M., Trtanj, J., Asher, J., Clay, M., Margolis, H. S., Hebbeler, A. M.,
39 George, D., & Chretien, J.-P. (2019). An open challenge to advance probabilistic forecasting for dengue epidemics.
40 *Proceedings of the National Academy of Sciences of the United States of America*, 116, 24268–24274. URL: [https:](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6883829/)
41 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6883829/>. doi:10.1073/pnas.1909865116.
42
43 Johnson, L. R., Gramacy, R. B., Cohen, J., Mordecai, E., Murdock, C., Rohr, J., Ryan, S. J., Stewart-
44 Ibarra, A. M., & Weikel, D. (2018). Phenomenological forecasting of disease incidence using het-
45 eroskedastic Gaussian processes: A dengue case study. *The Annals of Applied Statistics*, 12, 27–
46 66. URL: [https://projecteuclid.org/journals/annals-of-applied-statistics/volume-12/issue-1/](https://projecteuclid.org/journals/annals-of-applied-statistics/volume-12/issue-1/Phenomenological-forecasting-of-disease-incidence-using-heteroskedastic-Gaussian-processes/10.1214/17-AOAS1090.full)
47 [Phenomenological-forecasting-of-disease-incidence-using-heteroskedastic-Gaussian-processes/10.](https://projecteuclid.org/journals/annals-of-applied-statistics/volume-12/issue-1/Phenomenological-forecasting-of-disease-incidence-using-heteroskedastic-Gaussian-processes/10.1214/17-AOAS1090.full)
48 [1214/17-AOAS1090.full](https://projecteuclid.org/journals/annals-of-applied-statistics/volume-12/issue-1/Phenomenological-forecasting-of-disease-incidence-using-heteroskedastic-Gaussian-processes/10.1214/17-AOAS1090.full). doi:10.1214/17-AOAS1090. Publisher: Institute of Mathematical Statistics.
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4 Lauer, S. A., Sakrejda, K., Ray, E. L., Keegan, L. T., Bi, Q., Suangtho, P., Hinjoy, S., Iamsirithaworn, S., Suthachana,
5 S., Laosiritaworn, Y., Cummings, D. A., Lessler, J., & Reich, N. G. (2018). Prospective forecasts of annual
6 dengue hemorrhagic fever incidence in Thailand, 2010-2014. *Proceedings of the National Academy of Sciences*, *115*,
7 E2175–E2182. URL: <https://www.pnas.org/content/115/10/E2175>. doi:10.1073/pnas.1714457115. Publisher:
8 National Academy of Sciences _eprint: <https://www.pnas.org/content/115/10/E2175.full.pdf>.
9
10
11 Lega, J., & Brown, H. E. (2016). Data-driven outbreak forecasting with a simple nonlinear growth model. *Epidemics*,
12 *17*, 19–26. URL: <https://www.sciencedirect.com/science/article/pii/S1755436516300329>. doi:10.1016/j.
13 *epidem*.2016.10.002.
14
15 Lichtendahl Jr, K. C., Grushka-Cockayne, Y., & Winkler, R. L. (2013). Is it better to average probabilities or
16 quantiles? *Management Science*, *59*, 1594–1611.
17
18 Lipsitch, M., Finelli, L., Heffernan, R. T., Leung, G. M., & Redd; for the 2009 H1N1 Surveillance Group, S. C.
19 (2011). Improving the evidence base for decision making during a pandemic: The example of 2009 influenza
20 A/H1N1. *Biosecurity and Bioterrorism: Biodefense Strategy, Practice, and Science*, *9*, 89–115. URL: <https://www.liebertpub.com/doi/full/10.1089/bsp.2011.0007>. doi:10.1089/bsp.2011.0007.
21
22
23
24 Lowe, R., Coelho, C. A., Barcellos, C., Carvalho, M. S., Cato, R. D. C., Coelho, G. E., Ramalho, W. M., Bailey, T. C.,
25 Stephenson, D. B., & Rod, X. (2016). Evaluating probabilistic dengue risk forecasts from a prototype early warning
26 system for Brazil. *eLife*, *5*, e11285. URL: <https://doi.org/10.7554/eLife.11285>. doi:10.7554/eLife.11285.
27
28 McDonald, D. J., Bien, J., Green, A., Hu, A. J., DeFries, N., Hyun, S., Oliveira, N. L., Sharpnack, J., Tang,
29 J., Tibshirani, R., Ventura, V., Wasserman, L., & Tibshirani, R. J. (2021). Can auxiliary indicators improve
30 COVID-19 forecasting and hotspot prediction? *Proceedings of the National Academy of Sciences*, *118*. URL:
31 <https://www.pnas.org/content/118/51/e2111453118>. doi:10.1073/pnas.2111453118.
32
33
34 McGowan, C. J., Biggerstaff, M., Johansson, M., Apfeldorf, K. M., Ben-Nun, M., Brooks, L., Convertino, M.,
35 Erraguntla, M., Farrow, D. C., Freeze, J., Ghosh, S., Hyun, S., Kandula, S., Lega, J., Liu, Y., Michaud, N., Morita,
36 H., Niemi, J., Ramakrishnan, N., Ray, E. L., Reich, N. G., Riley, P., Shaman, J., Tibshirani, R., Vespignani,
37 A., Zhang, Q., & Reed, C. (2019). Collaborative efforts to forecast seasonal influenza in the United States,
38 2015-2016. *Scientific Reports*, *9*, 683. URL: <https://www.nature.com/articles/s41598-018-36361-9>. doi:10.
39 1038/s41598-018-36361-9.
40
41
42 Osthus, D., Gattiker, J., Priedhorsky, R., & Valle, S. Y. D. (2019). Dynamic Bayesian influenza fore-
43 casting in the United States with hierarchical discrepancy (with discussion). *Bayesian Analysis*,
44 *14*, 261–312. URL: [https://projecteuclid.org/journals/bayesian-analysis/volume-14/issue-1/
45 *Dynamic-Bayesian-Influenza-Forecasting-in-the-United-States-with-Hierarchical*/10.1214/
46 *18-BA1117.full*. doi:10.1214/18-BA1117. Publisher: International Society for Bayesian Analysis.
47
48
49 Osthus, D., Hickmann, K. S., Caragea, P. C., Higdon, D., & Del Valle, S. Y. \(2017\). Forecasting seasonal influenza
50 with a state-space SIR model. *The Annals of Applied Statistics*, *11*, 202–224. URL: \[https://www.ncbi.nlm.nih.
51 gov/pmc/articles/PMC5623938/\]\(https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5623938/\). doi:10.1214/16-AOAS1000.
52
53
54 Osthus, D., & Moran, K. R. \(2021\). Multiscale influenza forecasting. *Nature Communications*, *12*, 2991. URL:
55 <https://www.nature.com/articles/s41467-021-23234-5>. doi:10.1038/s41467-021-23234-5.
56
57 Pei, S., Kandula, S., Yang, W., & Shaman, J. \(2018\). Forecasting the spatial transmission of influenza in the
58
59
60
61
62
63
64
65](https://projecteuclid.org/journals/bayesian-analysis/volume-14/issue-1/Dynamic-Bayesian-Influenza-Forecasting-in-the-United-States-with-Hierarchical/10.1214/18-BA1117.full)

- 1
2
3
4 United States. *Proceedings of the National Academy of Sciences*, 115, 2752–2757. URL: <https://www.pnas.org/content/115/11/2752>. doi:10.1073/pnas.1708856115. Publisher: National Academy of Sciences Section:
5
6 Biological Sciences.
7
8 Polikar, R. (2006). Ensemble based systems in decision making. *IEEE Circuits and Systems Magazine*, 6, 21–45.
9 doi:10.1109/MCAS.2006.1688199. Conference Name: IEEE Circuits and Systems Magazine.
10
11 Pollett, S., Johansson, M. A., Reich, N. G., Brett-Major, D., Valle, S. Y. D., Venkatramanan, S., Lowe, R., Porco, T.,
12 Berry, I. M., Deshpande, A., Kraemer, M. U. G., Blazes, D. L., Pan-ngum, W., Vespigiani, A., Mate, S. E., Silal,
13 S. P., Kandula, S., Sippy, R., Quandelacy, T. M., Morgan, J. J., Ball, J., Morton, L. C., Althouse, B. M., Pavlin,
14 J., Panhuis, W. v., Riley, S., Biggerstaff, M., Viboud, C., Brady, O., & Rivers, C. (2021). Recommended report-
15 ing items for epidemic forecasting and prediction research: The EPIFORGE 2020 guidelines. *PLOS Medicine*,
16 18, e1003793. URL: <https://journals.plos.org/plosmedicine/article?id=10.1371/journal.pmed.1003793>.
17 doi:10.1371/journal.pmed.1003793.
18
19 Ranjan, R., & Gneiting, T. (2010). Combining probability forecasts. *Journal of the Royal*
20 *Statistical Society: Series B (Statistical Methodology)*, 72, 71–91. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1467-9868.2009.00726.x>. doi:10.1111/j.1467-9868.2009.00726.x. eprint:
21 <https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1467-9868.2009.00726.x>.
22
23 Ratcliff, R. (1979). Group reaction time distributions and an analysis of distribution statistics. *Psychological bulletin*,
24 86, 446–461.
25
26 Ray, E. (2020). reichlab/covidEnsembles: pre-publication release. URL: <https://zenodo.org/record/3963370>.
27 doi:10.5281/zenodo.3963370.
28
29 Ray, E. L. (2021). COVID-19 ensemble methods manuscript. URL: <https://zenodo.org/record/5784745>. doi:10.
30 5281/zenodo.5784745.
31
32 Ray, E. L., Brooks, L. C., Bien, J., Bracher, J., Gerding, A., Rumack, A., Biggerstaff,
33 M., Johansson, M. A., Tibshirani, R. J., & Reich, N. G. (2021 [Online]). Challenges
34 in training ensembles to forecast COVID-19 cases and deaths in the United States. In-
35 ternational Institute of Forecasters blog. URL: [https://forecasters.org/blog/2021/04/09/
36 challenges-in-training-ensembles-to-forecast-covid-19-cases-and-deaths-in-the-united-states/](https://forecasters.org/blog/2021/04/09/challenges-in-training-ensembles-to-forecast-covid-19-cases-and-deaths-in-the-united-states/).
37
38 Ray, E. L., & Reich, N. G. (2018). Prediction of infectious disease epidemics via weighted density ensembles.
39 *PLOS Computational Biology*, 14, 1–23. URL: <https://doi.org/10.1371/journal.pcbi.1005910>. doi:10.1371/
40 journal.pcbi.1005910.
41
42 Ray, E. L., Sakrejda, K., Lauer, S. A., Johansson, M. A., & Reich, N. G. (2017). Infectious disease prediction with
43 kernel conditional density estimation. *Statistics in medicine*, 36, 4908–4929. URL: [https://www.ncbi.nlm.nih.
44 gov/pmc/articles/PMC5771499/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5771499/). doi:10.1002/sim.7488.
45
46 Reich, N. G., Lauer, S. A., Sakrejda, K., Iamsirithaworn, S., Hinjoy, S., Suangtho, P., Suthachana, S., Clapham, H. E.,
47 Salje, H., Cummings, D. A. T., & Lessler, J. (2016). Challenges in Real-Time Prediction of Infectious Disease: A
48 Case Study of Dengue in Thailand. *PLOS Neglected Tropical Diseases*, 10, e0004761. URL: [https://journals.
49 plos.org/plosntds/article?id=10.1371/journal.pntd.0004761](https://journals.plos.org/plosntds/article?id=10.1371/journal.pntd.0004761). doi:10.1371/journal.pntd.0004761. Pub-
50 lisher: Public Library of Science.
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4 Reich, N. G., McGowan, C. J., Yamana, T. K., Tushar, A., Ray, E. L., Osthus, D., Kandula, S., Brooks,
5 L. C., Crawford-Crudell, W., Gibson, G. C., Moore, E., Silva, R., Biggerstaff, M., Johansson, M. A., Rosen-
6 feld, R., & Shaman, J. (2019). Accuracy of real-time multi-model ensemble forecasts for seasonal influenza in
7 the U.S. *PLOS Computational Biology*, *15*, e1007486. URL: [https://www.ncbi.nlm.nih.gov/pmc/articles/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6897420/)
8 [PMC6897420/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6897420/). doi:10.1371/journal.pcbi.1007486.
9
10 Reich, N. G., Tibshirani, R. J., Ray, E. L., & Rosenfeld, R. (2021 [Online]). On the predictability of
11 COVID-19. International Institute of Forecasters blog. URL: [https://forecasters.org/blog/2021/09/28/](https://forecasters.org/blog/2021/09/28/on-the-predictability-of-covid-19/)
12 [on-the-predictability-of-covid-19/](https://forecasters.org/blog/2021/09/28/on-the-predictability-of-covid-19/).
13
14 Reis, J., Yamana, T., Kandula, S., & Shaman, J. (2019). Superensemble forecast of respiratory syncytial virus
15 outbreaks at national, regional, and state levels in the United States. *Epidemics*, *26*, 1–8. URL: [https://www.](https://www.sciencedirect.com/science/article/pii/S1755436517301743)
16 [sciencedirect.com/science/article/pii/S1755436517301743](https://www.sciencedirect.com/science/article/pii/S1755436517301743). doi:10.1016/j.epidem.2018.07.001.
17
18 Shaman, J., & Karspeck, A. (2012). Forecasting seasonal outbreaks of influenza. *Proceedings of the National*
19 *Academy of Sciences*, *109*, 20425–20430. URL: <https://www.pnas.org/content/109/50/20425>. doi:10.1073/
20 [pnas.1208772109](https://www.pnas.org/content/109/50/20425). Publisher: National Academy of Sciences Section: Biological Sciences.
21
22 Taylor, J. W., & Taylor, K. S. (2021). Combining probabilistic forecasts of COVID-19 mortality in the United States.
23 *European Journal of Operational Research*, . doi:10.1016/j.ejor.2021.06.044.
24
25 Turtle, J., Riley, P., Ben-Nun, M., & Riley, S. (2021). Accurate influenza forecasts using type-specific incidence data
26 for small geographic units. *PLOS Computational Biology*, *17*, e1009230. URL: [https://journals.plos.org/](https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1009230)
27 [ploscompbiol/article?id=10.1371/journal.pcbi.1009230](https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1009230). doi:10.1371/journal.pcbi.1009230. Publisher:
28 Public Library of Science.
29
30 Viboud, C., Sun, K., Gaffey, R., Ajelli, M., Fumanelli, L., Merler, S., Zhang, Q., Chowell, G., Simonsen, L., &
31 Vespignani, A. (2018). The RAPIDD Ebola forecasting challenge: Synthesis and lessons learnt. *Epidemics*,
32 *22*, 13–21. URL: <https://www.sciencedirect.com/science/article/pii/S1755436517301275>. doi:10.1016/j.
33 [epidem.2017.08.002](https://www.sciencedirect.com/science/article/pii/S1755436517301275).
34
35 Vincent, S. B. (1912). *The Functions of the Vibrissae in the Behavior of the White Rat* volume 1. University of
36 Chicago.
37
38 Wallinga, J., Boven, M. v., & Lipsitch, M. (2010). Optimizing infectious disease interventions during an emerging
39 epidemic. *Proceedings of the National Academy of Sciences*, *107*, 923–928. URL: [https://www.pnas.org/content/](https://www.pnas.org/content/107/2/923)
40 [107/2/923](https://www.pnas.org/content/107/2/923). doi:10.1073/pnas.0908491107.
41
42 Yamana, T. K., Kandula, S., & Shaman, J. (2016). Superensemble forecasts of dengue outbreaks. *Journal of The*
43 *Royal Society Interface*, *13*, 20160410. URL: [https://royalsocietypublishing.org/doi/full/10.1098/rsif.](https://royalsocietypublishing.org/doi/full/10.1098/rsif.2016.0410)
44 [2016.0410](https://royalsocietypublishing.org/doi/full/10.1098/rsif.2016.0410). doi:10.1098/rsif.2016.0410. Publisher: Royal Society.
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Journal Pre-proof