

Nowcasting Reported COVID-19 Hospitalizations Using De-Identified, Aggregated Medical Insurance Claims Data

Xueda Shen¹ Aaron Rumack² Bryan Wilder² Ryan J. Tibshirani¹

¹University of California, Berkeley ²Carnegie Mellon University

Abstract

We propose, implement, and evaluate a method for nowcasting the daily number of new COVID-19 hospitalizations, at the level of individual US states, based on de-identified, aggregated medical insurance claims data. Our analysis proceeds under a hypothetical scenario in which, during the Delta wave, states only report data on the first day of each month, and on this day, report COVID-19 hospitalization counts for each day in the previous month. In this hypothetical scenario (just as in reality), medical insurance claims data continues to be available daily. At the beginning of each month, we train a regression model, using all data available thus far, to predict hospitalization counts from medical insurance claims. We then use this model to nowcast the (unseen) values of COVID-19 hospitalization counts from medical insurance claims, at each day in the following month. Our analysis uses properly-versioned data, which would have been available in real-time, at the time predictions are produced. In spite of the difficulties inherent to real-time estimation (e.g., latency and backfill) and the complex dynamics behind COVID-19 hospitalizations themselves, we find overall that medical insurance claims can be an accurate predictor of hospitalization reports, with mean absolute errors typically around 0.4 hospitalizations per 100,000 people, i.e., proportion of variance explained around 75%. Perhaps more importantly, we find that nowcasts made using medical insurance claims can qualitatively capture the dynamics (upswings and downswings) of hospitalization waves, which are key features that inform public health decision-making.

1 Introduction

Timely access to public health data is critical to enable informed decision making during infectious disease outbreaks. However, setting up and maintaining public health reporting pipelines can be a burden on the health system itself. For example, beginning in May 2020, hospitals in the US have been required to report data on COVID-19 hospitalizations to the Department of Health and Human Services (HHS) ([Department of Health and Human Services, 2023](#)). This data has been critical for understanding the state of the pandemic and the current load on the health system. But the frequent reporting required for up-to-date situational awareness (daily, throughout most of the pandemic) has been quite difficult to implement and maintain. In an effort to achieve compliance, the Centers for Medicare and Medicaid Services (CMS) issued regulations in August 2020 that threatened to expel hospitals from the Medicare program, and apply monetary penalties, if they failed to comply with daily COVID-19 reporting requirements. This was strongly and openly opposed by the American Hospital Association (AHA) ([Nickels, 2020](#)), but the regulations remained in place for nearly three years, lasting until the conclusion of the COVID-19 Public Health Emergency in May 2023.

A potential alternative lies in data streams which already exist and are already maintained for other purposes, and yet are relevant for inferring disease activity. One example is medical insurance claims, which are filed by healthcare providers to seek reimbursement from a patient’s insurance company for medical services performed. In this paper, we examine the use of *de-identified, aggregated* medical insurance claims data as a complement to public health reporting over the course of the pandemic. More specifically, we ask: if hospital reporting on COVID-19 would have been reduced in frequency from daily to monthly during the Delta wave, could we use signals derived from medical insurance claims to accurately *nowcast* COVID-19 hospitalization counts during the interim periods?

Medical insurance claims have long been utilized in public health policy analysis, ranging from economic implications of healthcare (e.g., recent examples include [Panczak et al. \(2018\)](#); [Li and Yang \(2021\)](#); [Sakai](#)

et al. (2019); Zheng and Peng (2021); Durizzo et al. (2022); Mori et al. (2022)), to examinations of treatment effectiveness and satisfaction (e.g., Nakayama et al. (2017); Jung et al. (2020); Geng et al. (2021); Yao et al. (2022); Song et al. (2023)). These works, however, are all *retrospective* in nature: they seek to develop an understanding of a particular phenomenon using data that would not have been available in real-time. In contrast, our goal to carry out an analysis that reflects *real-time* estimation, so that we can understand (to the best extent we can) how our models would perform if they were to be operationalized in the future, for true prospective nowcasting. This requires us to use properly-versioned data at all times in our analysis, that would have been available in real-time, at the time nowcasts are produced. This is true of all data sources in question, but it is especially crucial for medical insurance claims signals: these are subject to heavy revisions, as claims can be filed long after a service was performed (we describe this more concretely in what follows), altering previously-computed signal values. Therefore, using finalized (rather than properly-versioned) values of medical insurance claims signals for modeling and prediction can present a misleading picture of nowcasting performance. Figure 1 displays a clear example of this, using data from California.

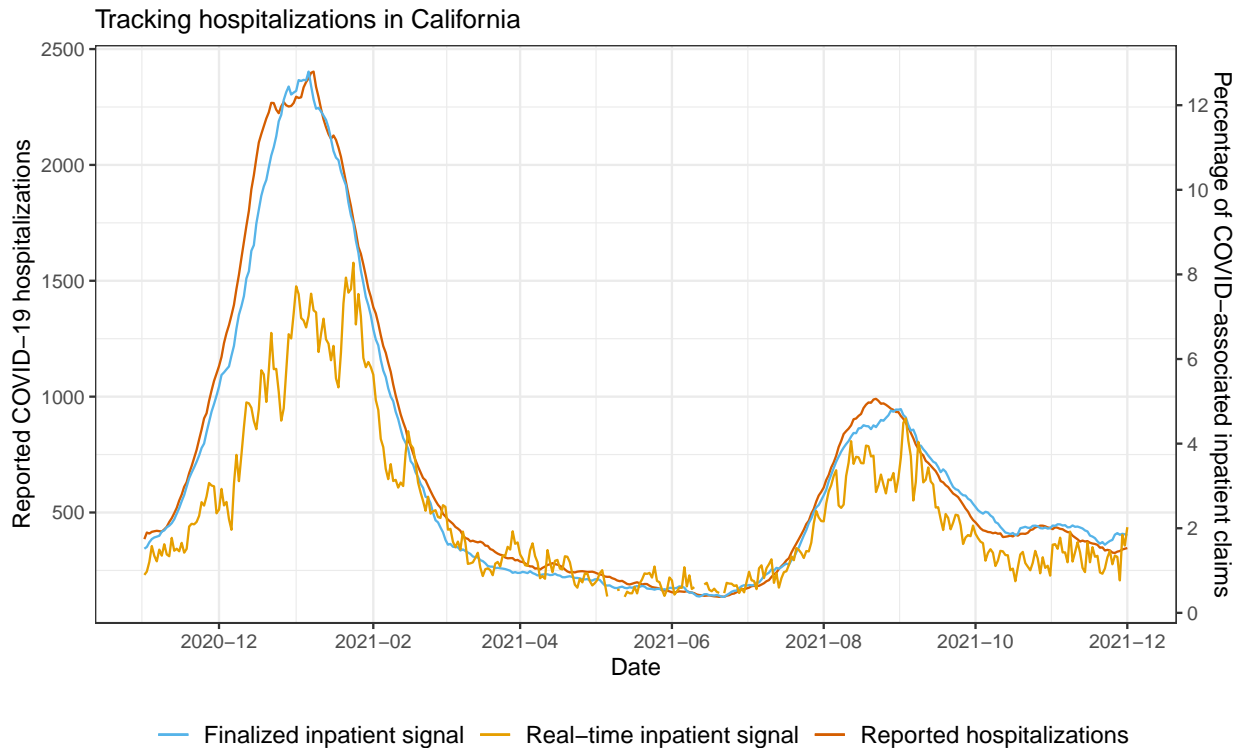


Figure 1: Reported daily COVID-19 hospitalizations, plotted alongside a signal derived from medical claims that measures daily COVID-associated inpatient admissions, for the state of California over a 1 year period during the pandemic. (All data is smoothed using a 7-day trailing average.) Two versions of the inpatient claims signal are shown: a real-time inpatient signal, which represents what would have been available at each reference date along the x-axis, and a finalized signal, which shows what would have been available 30 days after the reference date. We can see that the finalized signal tracks reported COVID-19 hospitalizations very well, but the concordance between reported hospitalizations and the real-time signal is much worse.

The importance of data versioning for epidemic tracking and prediction tasks is emphasized in Reinhart et al. (2021); McDonald et al. (2021), and the importance of leveraging existing healthcare data streams for epidemic surveillance, as a complement to traditional public health reporting, is motivated in Rosenfeld and Tibshirani (2021). Infectious disease nowcasting, using data from healthcare pipelines and also from a variety of other auxiliary data sources, has received increased attention over the last 10 years: see, e.g., Viboud et al. (2014); Smolinski et al. (2015); Farrow (2016); Santillana et al. (2016); Jahja et al. (2019); Yang et al. (2019); Ackley et al. (2020); Brooks (2020); Leuba et al. (2020); Radin et al. (2020), among others. Most work in this

area focuses on nowcasting as a means of producing up-to-date, high-resolution estimates of disease activity for a pathogen like influenza, where traditional public health surveillance streams are far from comprehensive. Our paper complements this line of work by studying COVID-19, where public health surveillance has been much more comprehensive (due to mandatory reporting). To our knowledge, we are the first to examine the effectiveness of nowcasting under a hypothetical scenario in which COVID-19 reporting *would have* been set up at a coarser frequency, which for hospitalization data, means less burden on hospitals themselves.

2 Methods

This section describes the data that we use, the hypothetical scenarios (for hospital reporting cadence) that we consider, and the nowcast and backcast models that we build and evaluate.

2.1 Data

Throughout, we restrict our attention to nowcasting state-level hospitalizations. This is because the medical claims signals that we describe below are not generally available at each US county. That said, in principle, the same ideas we describe in what follows could be applied to nowcast hospital reports in large counties (subject to enough claims data being available in order to form robust signals). We also restrict our attention to nowcasting daily reported hospitalizations between April 1, 2021 and August 1, 2023, though we use data back until November 1, 2020 for training models. We will revisit the choice of time periods shortly, when we describe the hypothetical scenarios, in the next subsection.

State-level, daily COVID-19 hospitalization reports are obtained from the HHS, accessed via the Delphi Epidata API (Farrow et al., 2015; Reinhart et al., 2021). We use $Y_{\ell,s}$ to denote the 7-day trailing average of *finalized* reported new COVID-19 hospitalization counts corresponding to location ℓ and time s . Averaging over 7-day trailing windows is mainly used as a smoother, and to account for weekday-weekend differences. Reported hospitalization counts are subject to revision (these are typically minor in comparison to revisions for claims signals), hence we also introduce notation to work with versioned data henceforth: we denote by $Y_{\ell,s}^{(t)}$ the 7-day trailing average of hospitalization counts for location ℓ and time s , but whose data version is *as of* time t . In this context, we often refer to s as the *reference date* and t the *issue date*. We use analogous notation and nomenclature for all versioned data in this paper.

De-identified, aggregated medical insurance claims data are provided to us by Change Healthcare. This data covers a sizeable fraction of healthcare providers in the US, enough for us to be able to generate robust signals of COVID-associated outpatient and inpatient activity for each U.S. state. Specifically, we consider the following two signals:

- $O_{\ell,s}$: the *finalized* percentage of outpatient claims in a 7-day trailing window with a confirmed COVID diagnostic code, corresponding to location ℓ and time s .
- $I_{\ell,s}$: the *finalized* percentage of inpatient claims in a 7-day trailing window with a COVID-associated diagnostic code, corresponding to location ℓ and time s .

As before, we use $O_{\ell,s}^{(t)}$ and $I_{\ell,s}^{(t)}$ to denote the *versioned* outpatient and inpatient claims signals, respectively, corresponding to issue date t . The use of versioned data is extremely important when working with claims-derived signals because medical insurance are quite often submitted and/or processed late, many days (and even months) after a given date of service. This process is referred to as *backfill*, and it generally has quite a pronounced effect on the outpatient and inpatient signals (note that the numerator and denominator defining these signals each get updated once new data comes in). We can look back at Figure 1 to see an example for the inpatient signal in California. What is labeled as the “real-time inpatient signal” in the figure is $I_{\ell,s}^{(s)}$ in our notation introduced here, for s ranging over the time values along the x-axis (and $\ell = \text{California}$). What is labeled as the “finalized inpatient signal” is actually $I_{\ell,s}^{(s+30)}$; we note that the choice of 30 days is arbitrary, made simply for visualization purposes, and it is not necessarily the case that all claims would have been filed after 30 days.

Sometimes backfill is so severe that no claims for reference date s are filed at all until a later issue date t , which we refer to as *latency*. If the latency is large enough, in particular, if no claims are available until a full

7 days after a given reference date, then the real-time claims signal value will be missing. This happens a few times in Figure 1 in May and June of 2021.

Like the HHS hospitalization reports, various signals derived from medical insurance claims are available in the Delphi Epidata API. The set of the diagnostic codes used to form the outpatient and inpatient signals are given in the API documentation: <https://cmu-delphi.github.io/delphi-epidata>; for convenience, we have relayed these definitions in Appendix A, and we have made all data (including properly-versioned data) used in our analysis available for download at: <https://github.com/cmu-delphi/hhs-nowcasting>.

2.2 Hypothetical scenarios

We consider two hypothetical scenarios in the analysis in what follows. Recall that our analysis examines nowcasting daily reported hospitalizations between April 1, 2021 and August 1, 2023, and we have training data all the way back until November 1, 2020.

Scenario 1: monthly updates. Our first hypothetical scenario spans nowcasting dates from April 1, 2021 to November 30, 2021. In this scenario, on the first day of each month during this period, we receive daily hospitalization counts for the previous month, and we nowcast and backcast the (unobserved) hospitalization counts for each following day in the given month, using the outpatient and inpatient claims signals described above. Note that the period in this scenario covers the Delta wave in the U.S. The start of this hypothetical scenario is chosen as April 1, 2021 so that we have enough data (reported hospitalizations and claims signals) to train our initial regression models for nowcasting. Recall that our data extends back through November 1, 2020, thus our initial training set includes the winter wave of 2020, which helps the initial models capture the relationship between reported hospitalization counts and claims signals.

To make this all more precise, let us introduce some notation. We drop reference to the location ℓ here and in what follows, whenever convenient (whenever it is not needed for the given explanation). Let t_0 be a date that marks the start of a month during the period of April 1, 2021 to November 30, 2021. Then on each day $t = t_0, t_0 + 1, \dots, t_1 - 1$, where t_1 marks the first day of the next month, we have access to:

- $\{Y_s^{(t_0)}\}_{s \leq t_0 - 1}$, hospitalization counts through day $t_0 - 1$, with versions as of day t_0 ; and
- $\{(I_s^{(t)}, O_s^{(t)})\}_{s \leq t}$, outpatient and inpatient signals through day t , with versions as of day t .

On each such day t , we use the data we have to train a regression model, call it f_t , to predict hospitalizations from claims signals. We use this to make nowcasts:

$$\hat{Y}_t^{(t)} = f_t(I_t^{(t)}, O_t^{(t)}),$$

and lag- k backcasts:

$$\hat{Y}_{t-k}^{(t)} = f_t(I_{t-k}^{(t)}, O_{t-k}^{(t)}),$$

for each $k = 1, \dots, 10$. (Clearly, a nowcast is equivalent to a lag-0 backcast.) This is repeated for each of the 9 months in the period spanned by this monthly-update scenario. The details of how regression models are trained and evaluated will be given in the next subsection.

Scenario 2: no updates. Our second hypothetical scenario spans nowcasting dates from December 1, 2021 to August 1, 2023. In this scenario, we receive reported hospitalizations up through November 30, 2021 with versions as of December 1, 2021, and we receive *no further* hospitalization counts after that. As before, we nowcast and backcast the (unobserved) hospitalization counts in the remaining period using the outpatient and inpatient claims signals. Note that the period in this scenario covers the Omicron wave in the U.S., and that this second scenario is generally far more challenging than the first scenario.

The data received and nowcasts and backcasts made in the no-update scenario can be written in precise notation, in fact, exactly as introduced in the description for the monthly-update scenario above, but now we fix t_0 at December 1, 2021, and make nowcasts and backcasts at each t from t_0 through the end of the period, August 31, 2023. The details of how regression models are trained and evaluated in this scenario will again be covered in the next subsection.

2.3 Regression model

As our basic working model, we predict hospitalizations using a linear combination of a set L^I of lags of the inpatient signal and a set L^O of lags of the outpatient signal,

$$f_t(I_s^{(t)}, O_s^{(t)}) = \beta_{t,0} + \sum_{j \in L^I} \beta_{t,j}^I I_{s-j}^{(t)} + \sum_{j \in L^O} \beta_{t,j}^O O_{s-j}^{(t)}. \quad (1)$$

where the coefficients $\beta_{t,0}$, $\beta_{t,j}^I$, and $\beta_{t,j}^O$ are estimated, i.e., the model f_t is fit, by training on historical data available at time t , either separately for each location ℓ (recall, the notational dependence on the location has been dropped for now), or in a way that pools data across locations. Details will be given below. First, we describe how we select the lag sets L^O , L^I for the inpatient and outpatient features.

2.3.1 Selecting feature lags

In order to select the lag sets L^O , L^I for the working model (1), we restrict our attention to training data *before* the first nowcast date of April 1, 2021 (so as to be careful not to overfit to the data in our nowcasting period, and reserve this for proper evaluation of our models in the two scenarios outlined previously).

To guide lag selection, we consider three metrics, which measure predictive power, feature stability, and data availability. Let $X_s^{(t)}$ denote a generic feature in our usual notation for versioned data, i.e., $X_s^{(t)} = I_s^{(t)}$ for the inpatient feature or $X_s^{(t)} = O_s^{(t)}$ for outpatient feature, we consider, as a function of lag j :

- correlation of $X_{t-j}^{(t)}$ with finalized hospitalizations Y_t , which is a measure of predictive power;
- correlation of $X_{t-j}^{(t)}$ with its own finalized value X_{t-j} , which is a measure of stability;
- the fraction of total claims used to compute the finalized signal X_{t-j} that are observed by time t .

Each metric is computed over time t (between November 1, 2020 and March 31, 2021), for each location, and the results are displayed in Figure 2, averaged over all locations (all states).

We can see that lag 6, for each of the inpatient and outpatient signals, provides a nice balance across the three metrics: maximum predictive power, and high stability and availability. We therefore include lag 6 in each of the sets L^I and L^O . To improve robustness and capture longer-range dependencies, we also include lags 13 and 20 in each of L^I and L^O , which is roughly consistent with choices of feature lags in other basic epidemic prediction models (e.g., McDonald et al. (2021)). The choice of 7-day spacing here is also motivated by the desire to limit correlations between features in (1); recall, the inpatient and outpatient signals are computed using a trailing 7-day window of claims data.

Lastly, we note that these lag sets have not really been directly optimized for maximum nowcasting and backcasting accuracy. We have avoided doing this for simplicity, and a more direct lag optimization over the training data (pre April 1, 2021) would likely only improve our predictive accuracy in general.

2.3.2 Training the regression model

The coefficients $\beta_{t,0}$, $\beta_{t,j}^I$, $\beta_{t,j}^O$ in (1) are fit by solving the following weighted least squares problem at time t :

$$\underset{\beta_{t,0}, \beta_{t,j}^I, \beta_{t,j}^O}{\text{minimize}} \quad \sum_{s < t_0} w_\gamma(t_0 - s) \cdot \left(Y_s^{(t_0)} - \beta_{t,0} - \sum_{j \in L^I} \beta_{t,j}^I I_{s-j}^{(t)} - \sum_{j \in L^O} \beta_{t,j}^O O_{s-j}^{(t)} \right)^2, \quad (2)$$

where we use exponentially decaying weights:

$$w_\gamma(u) = \exp(-\gamma u) \quad (3)$$

for a tuning parameter $\gamma \geq 0$. The index t_0 in (2) marks the latest observation boundary for hospitalization reports before time t : this is either the start of the month in scenario 1 (where we receive monthly updates), or December 1, 2021 in scenario 2 (where we receive no further updates). Hence, in other words, we fit the coefficients in (2) by minimizing the weighted mean squared error of our working regression model, over all dates at which response values (reported hospitalization counts) are available. Note carefully that in (2) we

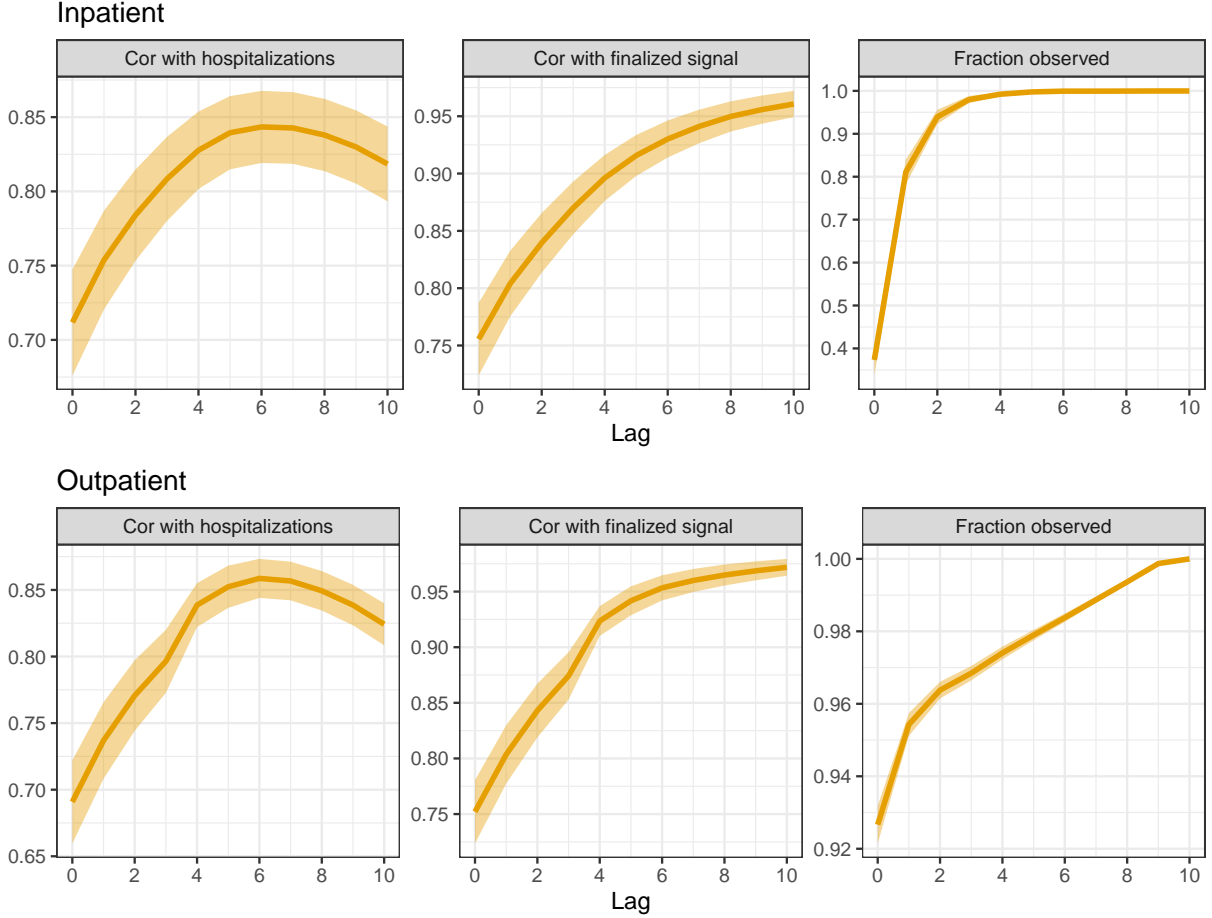


Figure 2: Analysis to help guide lag selection for inpatient and outpatient features in the working model (1). The rows correspond to different features—inpatient or outpatient, and the columns to different metrics, as explained precisely in the main text.

only use properly-versioned data that would have been available at t . If any such data, response or feature values, are missing at time t then the corresponding summand is simply omitted from (2).

As a default, we fit the model by solving (2) separately for each location ℓ (recall, we have suppressed the dependence on ℓ in the notation for simplicity). We call this the *state-level* model. Later, we discuss schemes for pooling training data across locations. Next, we focus on the decay parameter γ in (2), (3).

2.3.3 Selecting the decay parameter by cross-validation

Before describing how we select γ in the exponential weights (3) that are used in the weighted least squares problem (2), we pause to discuss the question: why is it useful to use decaying observation weights in the first place? The reason is because the the features—lagged versions of the inpatient and outpatient medical insurance claims signals, and response—reported hospitalizations, need not be jointly stationary, i.e., their relationship may be changing over time.

In fact, we can already see clear evidence of nonstationarity in the relationship between the inpatient signal (red) and reported hospitalizations (gray) in Figure 1. If we were to regress reported hospitalizations on this inpatient signal alone, then it looks like a regression coefficient of about 1 would be appropriate for the period of April–July 2021, but this would be far too small for the first hospitalization wave starting in December 2020 (and to a lesser extent, too small for the second wave starting in August 2021).

By allowing γ itself change over time, we can adapt to the degree of nonstationarity at any point in time, i.e., we can adapt to the amount of past training data that is relevant for the current prediction. Indeed, we

will select γ in a dynamic, time-varying fashion using cross-validation (CV). Given the sequential nature of our prediction problem, we use a version of CV that is purely forward-looking, and is sometimes called *time series cross-validation* in the literature (Hyndman and Athanasopoulos, 2021). This works as follows. Let t_0 denote the index that marks the most recent observation boundary for reported hospitalizations, and t_{-1}, t_{-2} denote the indices that mark the previous two observation boundaries before t_0 . Then, for each γ in a grid Γ of tuning parameter values, we carry out the following procedure:

- for each $t \in [t_{-2}, t_{-1})$:
 - fit a regression model by solving:

$$\underset{\beta_{t,0}, \beta_{t,j}^I, \beta_{t,j}^O}{\text{minimize}} \quad \sum_{s < t_{-2}} w_\gamma(t_{-2} - s) \cdot \left(Y_s^{(t_{-2})} - \beta_{t,0} - \sum_{j \in L^I} \beta_{t,j}^I I_{s-j}^{(t)} - \sum_{j \in L^O} \beta_{t,j}^O O_{s-j}^{(t)} \right)^2;$$

- produce backcasts $\hat{Y}_s^{(t)}$, $s \in [t_{-2}, t]$;

- for each $t \in [t_{-1}, t_0)$:
 - fit a regression model by solving:

$$\underset{\beta_{t,0}, \beta_{t,j}^I, \beta_{t,j}^O}{\text{minimize}} \quad \sum_{s < t_{-1}} w_\gamma(t_{-1} - s) \cdot \left(Y_s^{(t_{-1})} - \beta_{t,0} - \sum_{j \in L^I} \beta_{t,j}^I I_{s-j}^{(t)} - \sum_{j \in L^O} \beta_{t,j}^O O_{s-j}^{(t)} \right)^2;$$

- produce backcasts $\hat{Y}_s^{(t)}$, $s \in [t_{-1}, t]$;

- compute the mean absolute error (MAE) of all backcasts made at all times $t \in [t_{-2}, t_0)$, against reported hospitalization counts $Y_t^{(t_0)}$, $t \in [t_{-2}, t_0)$.

In other words, this procedure uses the last 2 months of data as a validation set for tuning γ . Ultimately, we choose $\gamma \in \Gamma$ that minimizes the MAE computed in the last step above. To be clear, after choosing γ in this way, we then fit the model by solving (2) and use this to make backcasts at each time $t \in [t_0, t_1)$, where t_1 is the observation boundary after t_0 .

In our experiments, we take the set Γ to be a grid of 25 evenly-spaced values between 0 and γ_{\max} , where the latter is defined as the value of γ such that the effective sample size of the weight sequence,

$$\frac{\|w_\gamma(t_0 - \cdot)\|_1^2}{\|w_\gamma(t_0 - \cdot)\|_2^2} = \frac{(\sum_{s < t_0} |w_\gamma(t_0 - s)|)^2}{\sum_{s < t_0} |w_\gamma(t_0 - s)|^2},$$

is equal to 30. This is a heuristic upper bound for the “reasonable” range of γ values, under the rationale that using at least (effectively) 1 month of training data helps to avoid fitted models that are too volatile.

2.3.4 Stabilizing predictions by geo-pooling

To borrow strength across locations, we consider fitting the regression model at each time t by pooling data across locations, which we refer to as the *geo-pooled* model. To be precise, instead of solving (2) per location ℓ , here we instead solve:

$$\underset{\beta_{t,0}, \beta_{t,j}^I, \beta_{t,j}^O}{\text{minimize}} \quad \sum_{\ell} \sum_{s < t_0} w_\gamma(t_0 - s) \cdot \left(\tilde{Y}_{\ell,s}^{(t_0)} - \beta_{t,0} - \sum_{j \in L^I} \beta_{t,j}^I I_{\ell,s-j}^{(t)} + \sum_{j \in L^O} \beta_{t,j}^O O_{\ell,s-j}^{(t)} \right)^2. \quad (4)$$

Note that in (4), the training set includes data from all locations ℓ . Importantly, in (4), the response $\tilde{Y}_{\ell,s}^{(t)}$ is the *hospitalization rate* in location ℓ at reference date s , as of time t , which is defined as the hospitalization count per 100,000 people, i.e.,

$$\tilde{Y}_{\ell,s}^{(t)} = 10^5 \cdot \frac{Y_{\ell,s}^{(t)}}{N_\ell},$$

where N_ℓ is the population of location ℓ . The use of rates rather than counts in the pooled regression (4) is critical, because otherwise the coefficients would have different meanings in different locations, and pooling across locations would not make sense. After solving (4), in order to use the fitted model to make predictions of hospitalization counts at a given location ℓ , we would then need to rescale the predictions by $N_\ell/10^5$.

Finally, we also consider a *mixed* model, which linearly combines the state-level $\hat{Y}_{\ell,s}^{(t),\text{state}}$ and geo-pooled $\hat{Y}_{\ell,s}^{(t),\text{pooled}}$ predictions, via

$$\hat{Y}_{\ell,s}^{(t),\text{mixed}} = \alpha \hat{Y}_{\ell,s}^{(t),\text{state}} + (1 - \alpha) \hat{Y}_{\ell,s}^{(t),\text{pooled}}, \quad (5)$$

for a mixing parameter $\alpha \in [0, 1]$. We select α , separately per state, using the same cross-validation strategy described previously, where we tune α over a grid of 50 evenly-spaced values between 0 and 1. (We tune over γ separately for each of the state-level and geo-pooled models, and then tune over α .)

3 Results

This section examines our backcasting results in scenario 1 (monthly-update) and scenario 2 (no-update), through quantitative backcast error analysis and qualitative inspection of the nowcast dynamics over time.

3.1 Scenario 1: backcast error analysis

Figure 3 shows the MAE of all backcasts, as a function of lag $k = 0, \dots, 10$, from the state-level, geo-pooled, and mixed models over the monthly-update period (which, recall spans April 1, 2021 to November 30, 2021). To be clear, for each model, this is computed by averaging the absolute error of backcasts made to *finalized* hospitalization counts, per state. The left panel displays the average of these state MAE values, and as well as standard error bands. The right panel is similar but normalizes the state MAE values by state population times 10^5 , thus considering MAE on the scale of hospitalization rates (rather than counts).

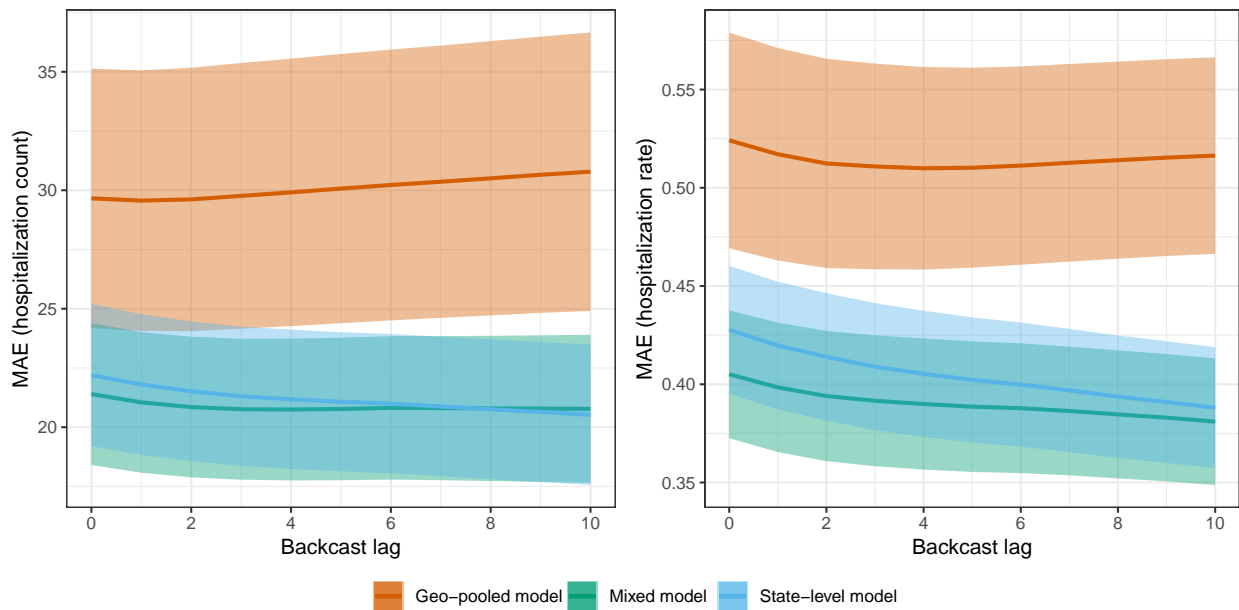


Figure 3: MAE as a function of backcast lag for the state-level, geo-pooled, and mixed models in scenario 1, the monthly-update period. The shaded regions show ± 1 standard error bands, over the state MAE values.

The geo-pooled model is clearly worse in terms of MAE than the mixed and state-level models. On the counts scale (left panel), the state-level model has slightly better MAE than the mixed model; on the rates scale (right panel), the opposite is true. This is because the mixed model generally performs better for smaller states, where predictions tend to be more volatile and shrinking towards the geo-pooled model helps reduce

variance. This is not reflected in the MAE plot on the left panel, since poor backcast performance on small states contributes little when measured in terms of counts.

In absolute terms, the backcasts made by the state-level and mixed models are generally quite accurate. For example, the nowcasts (backcasts at lag 0) made by the state-level and mixed models have MAEs of 0.428 and 0.405, respectively, on the scale of hospitalization rates. These correspond to proportions of variance explained (PVEs) of 71.6% and 75.4%, respectively. Even the geo-pooled model, which is relatively quite a bit worse, produces nowcasts with an MAE of 0.524 on the scale of hospitalization rates, which corresponds to a still a respectable PVE of 61.1%.

3.2 Scenario 1: illustrative nowcast examples

We examine nowcasts made by the state-level and mixed models during the monthly-update period, in four states: California (CA), Kentucky (KY), Vermont (VT), and New York (NY). The first three are chosen to demonstrate the qualitative behavior of nowcasts in states of different sizes, with CA having the largest population, KY having roughly median population, and VT having the smallest. NY is chosen because it represents somewhat of a failure case for the robustness of nowcasts from the state-level model.

California, Kentucky, Vermont. Figure 4 displays state-level and mixed model nowcasts for CA, KY, and VT. To be clear, in each panel of the figure, the nowcasts use *real-time* inpatient and claims signals as predictive features, and the models behind these nowcasts are trained using hospitalization data up to the latest observation boundary (marked by dotted vertical lines). The shaded bars in the figure display *finalized* reported hospitalizations, the target of ultimate interest.

In CA and KY, where the Delta wave is prominent (roughly July–November), we can see that both the state-level and mixed model nowcasts track the dynamics of the Delta wave. For example, looking at July 1, 2021 in CA, hospitalization reports from the previous months (which is all that would have been available at that time) show no indication of an upswing to come. Still, the nowcasts during July present a clear upward trend, which means that policy-makers would know from these nowcasts that a wave is underway. As we see it, evaluating nowcasts for their qualitative shape (in comparison to hospitalization waves) is important—just as much as (if not more so than) numeric error analysis, as in the last subsection.

Being much smaller, VT has on the order of ten reported COVID-19 hospitalizations daily, not several hundreds, as in KY or CA. Accordingly, we can see in Figure 4 that the (finalized) reported hospitalizations curve in VT is much noisier—we do not really see a clear Delta wave (instead, an increasing but quite noisy trend from August through December). The claims signals are also more noisy in a small state like VT, since they are ratios of small counts. Therefore, both the response and features are more volatile in the prediction problem for VT, and not surprisingly, the nowcasts here appear qualitatively worse. That said, the nowcasts still roughly track the gross trends in reported hospitalizations: higher in April and May, lower in June and July, and growing again in August onward.

Generally, the discussion above translates to the rest of the US: nowcasts tend to be in good qualitative agreement with hospitalization waves in larger states, and less so in smaller states. Appendix C provides a full set of nowcast and backcast plots in all 50 states, from the mixed model, under scenario 1. Occasionally, we find that the trend in the predictions disagrees with that in hospitalization reports in a given month, but this gets corrected at the next observation boundary, once new training data is available and the regression model is refit. We also find that backcasts at lag 5 or 10 tend to be smoother than nowcasts.

New York. Figure 4 displays state-level and mixed model nowcasts for NY. This is a notable example of a failure in robustness of the state-level model: its nowcasts for a good part of the month of June are actually *negative*, and are truncated at zero for visualization purposes. This is due to a large and systematic revision that occurred on June 8 for the outpatient signal, where this signal’s values for reference dates in May were revised upward. Moreover, when the state-level model is fit to data through the end of May, the regression coefficient on the largest outpatient feature lag ends up being negative, and consequently the upward revision of past feature values on June 8 introduces a strong downward bias in the subsequent predictions. Supporting analysis for this explanation is given in Appendix B.

As we can see in the figure, the mixed model produces nowcasts that are much more stable in the month of June. However, the mixed model and even the geo-pooled model are themselves not immune to large and

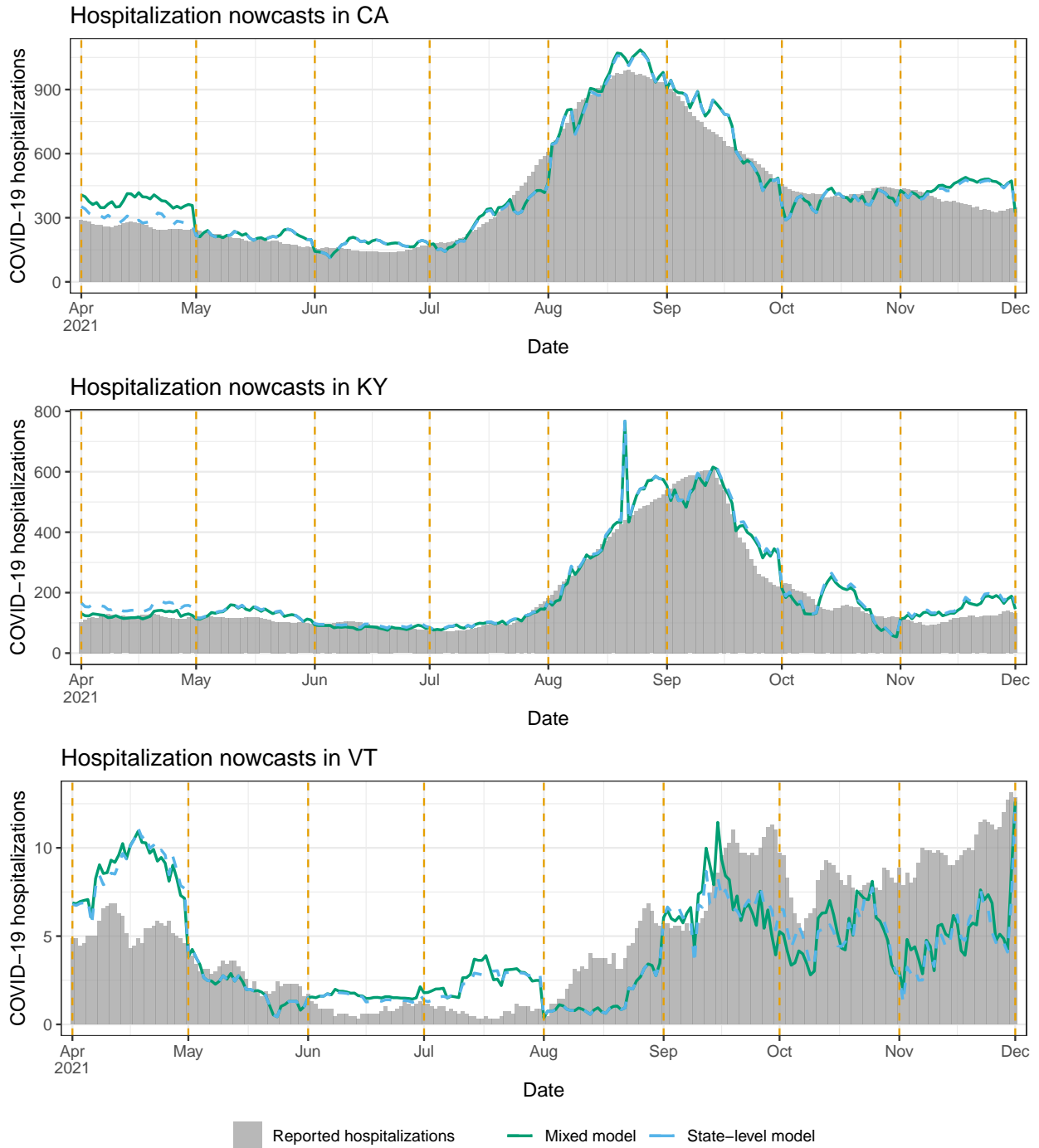


Figure 4: Nowcasts from the state-level and mixed models in scenario 1, the monthly-update period, for CA, KY, and VT. The dotted vertical lines mark the observation boundaries—when hospitalization reports for the previous month are received, and models are retrained.

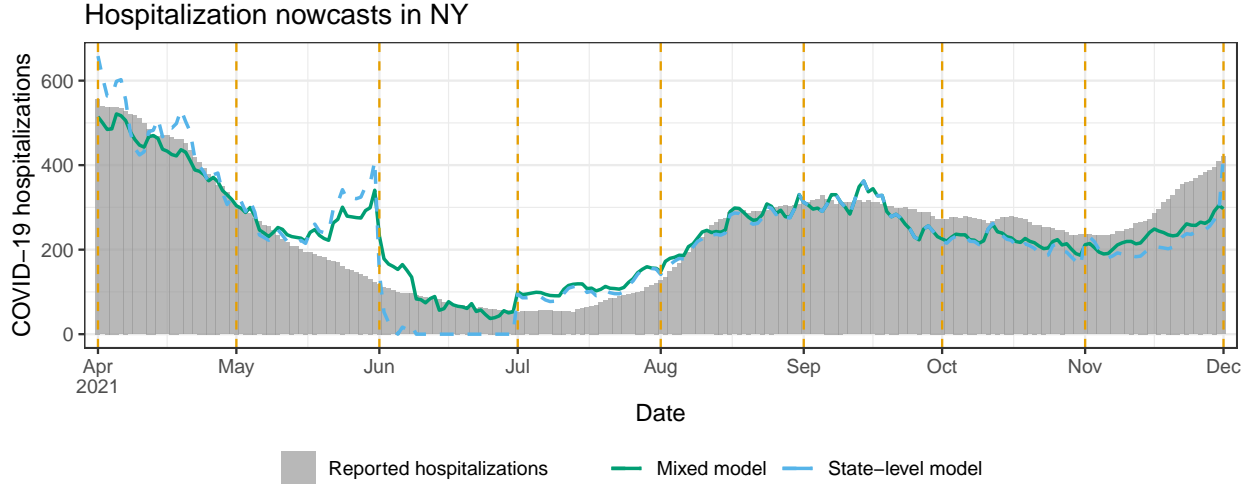


Figure 5: Nowcasts from the state-level and mixed models in scenario 1, the monthly-update period, for NY. The dotted vertical lines mark the observation boundaries, as in Figure 4.

erratic revisions in the input features at prediction time. These revisions can result in erratic nowcasts and backcasts. To mitigate this, we could turn to models that fuse information across a wider variety of auxiliary signals (where ideally, some of these signals would have less severe backfill than claims-based signals), an idea that we return to in the discussion section.

3.3 Scenario 2: backcast error analysis

Figure 6 shows the MAE of all backcasts, as a function of lag $k = 0, \dots, 10$, from the state-level, geo-pooled, and mixed models over the no-update period (spanning December 1, 2021 to August 31, 2023). The format is as in Figure 3. We see two notable differences compared to the results in the monthly-update period. First, as expected, each model performs worse than it did in Figure 6. On the scale of hospitalization rates (right panel), the range of MAEs has jumped from (roughly) 0.4–0.53 in the monthly-update period to 0.62–0.72 in the no-update period, and correspondingly, the PVEs have dropped from (roughly) 61–75% to 10–37%.

Second (and perhaps a bit more surprising), we see in the no-update period that the state-level model performs clearly worse than the geo-pooled and mixed models when MAE is measured either on the counts or rates scale. The mixed model MAE is somewhat better than the geo-pooled model MAE with respect to counts, while the two are basically the same with respect to rates. It is encouraging to see that the mixed model performs competitively across both monthly-update and no-update scenarios: this model—equipped with the CV procedure for tuning α —is able to effectively adapt the strengths of local and global training approaches (underlying the state-level and geo-pooled models) to the task at hand, in order to yield accurate predictions in an average-case sense.

3.4 Scenario 2: illustrative nowcast examples

We examine nowcasts made by the geo-pooled and mixed models for CA, KY, and VT. NY, which served as failure case in the monthly-update period, is not shown here, because KY itself provides such an example: as we will see, the mixed model lacks robustness over a part of the Omicron wave (meanwhile, the mixed model performs fine for NY throughout the no-update period—as shown in the appendix).

Figure 7 displays the nowcasts for CA, KY, and VT, with the same general format as in Figure 4. CA is a clear success case: both geo-pooled and mixed models capture the dynamics of the Omicron wave faithfully, and also track the summer 2023 and winter 2023 waves, despite (recall) receiving no reported hospitalizations past December 1, 2022. VT, as in the monthly-update period, is a challenging case because it corresponds to much noisier prediction problem, and the nowcasts here look qualitatively worse overall. However, the mixed model nowcasts still pick up the Omicron wave. KY represents a failure case: the mixed model nowcasts for

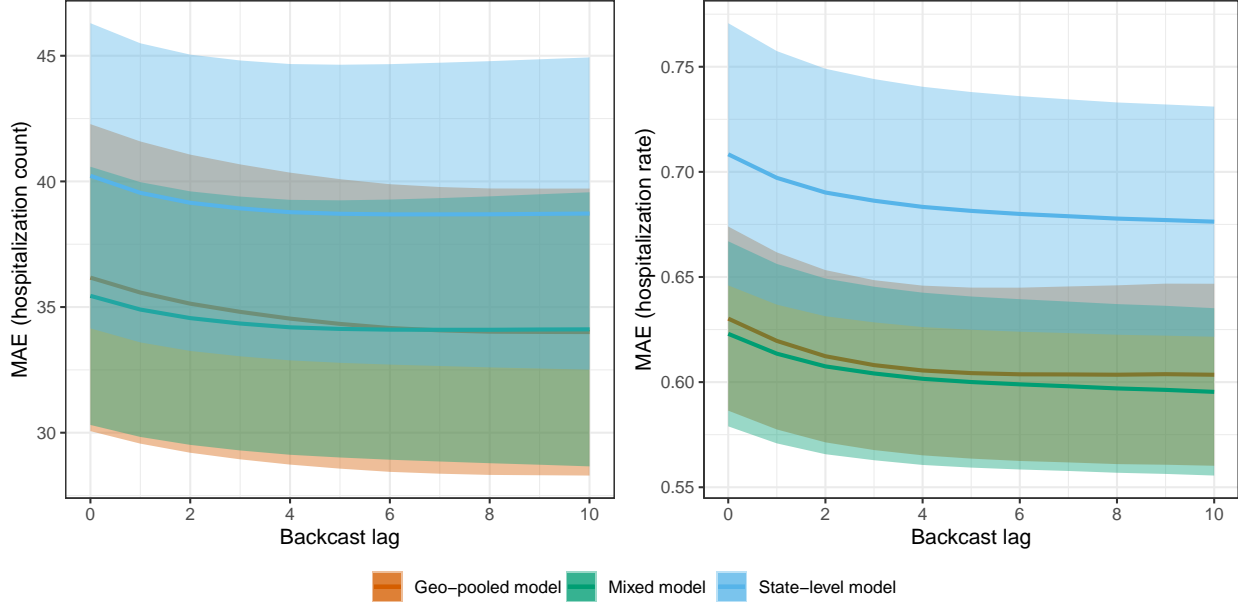


Figure 6: MAE as a function of backcast lag for the state-level, geo-pooled, and mixed models in scenario 2, the no-update period. The shaded regions show ± 1 standard error bands, over the state MAE values.

all of February are *negative* and truncated at zero for the visualization. What this is actually demonstrating is a failure of the state-level model, and simultaneously, a failure of tuning of the mixing parameter α . The mixed model here ends up placing a large weight on state-level predictions (not shown), which are themselves volatile for reasons similar to what happens in NY during the monthly-update period—large revisions to the input features at prediction time.

The volatility of the mixed model nowcasts in KY during Omicron provides an important perspective: in the MAE sense, we found in Figure 6 that the mixed model (with its CV tuning for α) successfully navigates the strengths and weaknesses of the geo-pooled and state-level models; and yet for specific states and periods of time, the mixed model can still lack robustness. To be clear, such fragility is not limited to the no-update period, and it could have happened in NY in the monthly-update period, had the CV tuning procedure for α not downweighted the state-level predictions so heavily. Improving robustness is a direction for future work, and we revisit this topic in the discussion section.

Lastly, in Appendix D, we again provide a full set of nowcast and backcast plots in all 50 states, from the mixed model, under scenario 2. For larger states, we find that the nowcasts generally trace out the Omicron wave, but for smaller states, this happens less consistently and the nowcasts look more noisy. Backcasts at lag 5 or 10 tend to smooth out the nowcasts, though not dramatically.

3.5 Ablation study

To examine the importance of some of our modeling choices (as described in the methods section), we carry out an ablation study in which we remove a particular component of the model, use a simpler alternative in its place, and evaluate the result in terms of MAE. In particular, we consider four ablated models:

1. Unweighted, all past: we fit the regression model (2) without observation weights, using all past reported hospitalization data available (considering summands $s < t$ in (2)).
2. Unweighted, two months: we fit the regression model (2) without observation weights, using reported hospitalization data available for the latest two months (considering summands $s \in [t-2, t_0]$ in (2)).
3. Weighted, inpatient only: we fit the regression model (2) using only the lags of the inpatient signal.
4. Weighted, outpatient only: we fit the regression model (2) using only the lags of the outpatient signal.

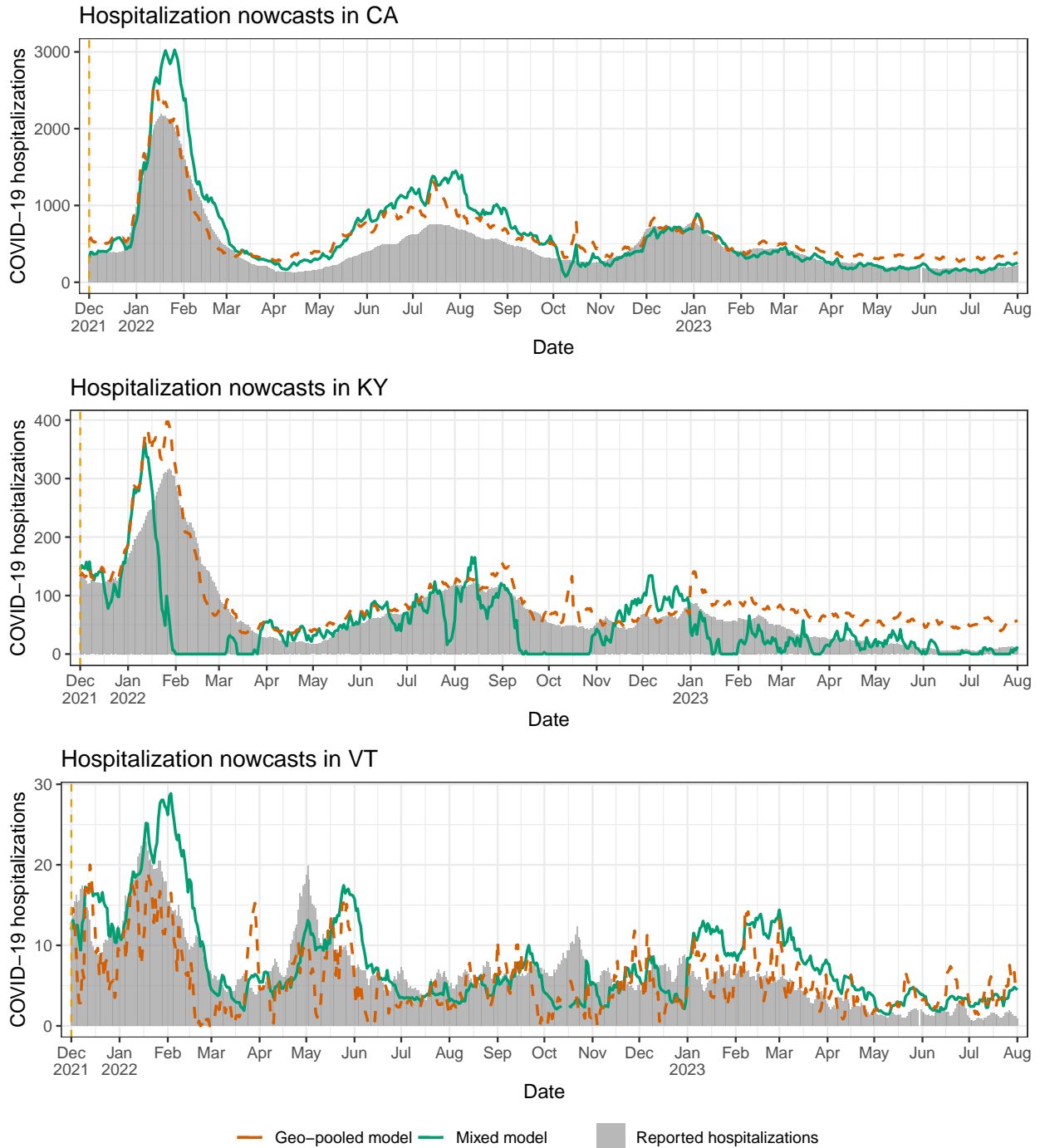


Figure 7: Nowcasts from the geo-pooled and mixed models in scenario 2, the no-update period, for CA, KY, and VT. The dotted vertical line (on the left side of the plot) marks the sole observation boundary in this period—hospitalization reports for all dates prior to this are available, but no reports are received after that.

Table 1 compares the performance of the state-level model to these ablated models, reporting the 25%, 50%, 75% quantiles of MAE values, over all backcast lags $k = 0, \dots, 10$ and all states. The state-level model, which uses decaying weights (with CV tuning for the decay parameter), and inpatient and outpatient features together, performs the best throughout. This emphasizes the importance of each of these two aspects of the original model design.

MAE (hospitalization count)				MAE (hospitalization rate)			
Model	25%	50%	75%	Model	25%	50%	75%
All-past	5.00	14.25	37.2	All-past	0.177	0.386	0.740
Two-month	9.05	27.4	66.7	Two-month	0.280	0.731	1.736
State-level	4.12	12.0	32.5	State-level	0.135	0.325	0.680
Inpat-only	4.91	13.8	33.1	Inpat-only	0.154	0.346	0.725
Outpat-only	5.46	16.3	41.0	Outpat-only	0.171	0.413	0.905

Table 1: MAE quartiles from the state-level and ablated models, pooled over all backcast lags and all states.

4 Discussion

This paper demonstrates that relatively simple regression models and medical insurance claims data can be used to provide accurate real-time estimates of reported COVID-19 hospitalizations, in different hypothetical scenarios in which hospitalization reporting is dramatically reduced in frequency or shut down completely. In general, leveraging statistical models that use auxiliary data for nowcasting, so we may then reduce reporting frequency and hence reduce the burden that reporting entails, may be a favorable tradeoff for public health agencies to consider.

Medical insurance claims are certainly not the only relevant auxiliary data stream for tracking COVID-19 hospitalizations, and this analysis may be repeated with a number of other auxiliary signals. In operational systems in public health, robustness is arguably more important than average-case performance, and failure examples of the proposed nowcasting methods, as seen in the results section (NY in scenario 1, and KY in scenario 2) would likely be concerning to public health decision-makers. These failure examples were driven by large revisions to the claims signals that occurred at certain points in time. Any model that makes linear predictions from a given set of features can suffer from erratic behavior if these features are (possibly) subject to large fluctuations at prediction time. Combining multiple nowcasts built from different auxiliary signals is a way to improve robustness, especially when some of these auxiliary signals are more stable and less subject to heavy revisions. Signals derived from electronic medical records (EMR), medical device data, and internet search queries all typically have less backfill compare to medical insurance claims signals.

Model combination methods—which can go by different names: aggregation, ensembling, or fusion—have been quite successful in influenza nowcasting systems in pre-pandemic years (e.g., [Farrow \(2016\)](#); [Jahja et al. \(2019\)](#)). Ensembles of COVID-19 forecasts have likewise demonstrated considerable robustness in comparison to the constituent forecasts (e.g., [Cramer et al. \(2022\)](#); [Ray et al. \(2023\)](#)). An important direction for future work is to incorporate similar ideas into the nowcasting settings considered in this paper, in an effort to move towards more reliable and robust systems.

Acknowledgements

We would like to thank members of the Delphi research group for valuable feedback, and Change Healthcare and Optum/United Health Group for their invaluable data partnership and collaboration. This work was supported by Centers for Disease Control and Prevention (CDC) grant no. 75D30123C15907.

References

- A. F. Ackley, S. Pilewski, V. S. Petrovic, L. Worden, E. Murray, and T. C. Porco. Assessing the utility of a smart thermometer and mobile application as a surveillance tool for influenza and influenza-like illness. *Health Informatics Journal*, 26(3):2148–2158, 2020.

- L. C. Brooks. *Pancasting: Forecasting epidemics from provisional data*. PhD thesis, Carnegie Mellon University, 2020.
- E. Y. Cramer, E. L. Ray, V. K. Lopez, J. Bracher, A. Brennen, A. J. C. Rivadeneira, A. Gerding, T. Gneiting, K. H. House, Y. Huang, D. Jayawardena, A. H. Kanji, A. Khandelwal, K. Le, A. Mühlemann, J. Niemi, A. Shah, A. Stark, Y. Wang, N. Wattanachit, M. W. Zorn, Y. Gu, S. Jain, N. Bannur, A. Deva, M. Kulkarni, S. Merugu, A. Raval, S. Shingi, A. Tiwari, J. White, S. Woody, M. Dahan, S. Fox, K. Gaither, M. Lachmann, L. A. Meyers, J. G. Scott, M. Tec, A. Srivastava, G. E. George, J. C. Cegan, I. D. Dettwiler, W. P. England, M. W. Farthing, R. H. Hunter, B. Lafferty, I. Linkov, M. L. Mayo, M. D. Parno, M. A. Rowland, B. D. Trump, S. M. Corsetti, T. M. Baer, M. C. Eisenberg, K. Falb, Y. Huang, E. T. Martin, E. McCauley, R. L. Myers, T. Schwarz, D. Sheldon, G. C. Gibson, R. Yu, L. Gao, Y. Ma, D. Wu, X. Yan, X. Jin, Y.-X. Wang, Y. Chen, L. Guo, Y. Zhao, Q. Gu, J. Chen, L. Wang, P. Xu, W. Zhang, D. Zou, H. Biegel, J. Lega, T. L. Snyder, D. D. Wilson, S. McConnell, R. Walraven, Y. Shi, X. Ban, Q.-J. Hong, S. Kong, J. A. Turtle, M. Ben-Nun, P. Riley, S. Riley, U. Koyluoglu, D. DesRoches, B. Hamory, C. Kyriakides, H. Leis, J. Milliken, M. Moloney, J. Morgan, G. Ozcan, C. Schrader, E. Shakhnovich, D. Siegel, R. Spatz, C. Stiefeling, B. Wilkinson, A. Wong, Z. Gao, J. Bian, W. Cao, J. L. Ferres, C. Li, T.-Y. Liu, X. Xie, S. Zhang, S. Zheng, A. Vespignani, M. Chinazzi, J. T. Davis, K. Mu, A. P. y. Piontti, X. Xiong, A. Zheng, J. Baek, V. Farias, A. Georgescu, R. Levi, D. Sinha, J. Wilde, N. D. Penna, L. A. Celi, S. Sundar, S. Cavany, G. España, S. Moore, R. Oidtman, A. Perkins, D. Osthus, L. Castro, G. Fairchild, I. Michaud, D. Karlen, E. C. Lee, J. Dent, K. H. Grantz, J. Kaminsky, K. Kaminsky, L. T. Keegan, S. A. Lauer, J. C. Lemaitre, J. Lessler, H. R. Meredith, J. Perez-Saez, S. Shah, C. P. Smith, S. A. Truelove, J. Wills, M. Kinsey, R. Obrecht, K. Tallaksen, J. C. Burant, L. Wang, L. Gao, Z. Gu, M. Kim, X. Li, G. Wang, Y. Wang, S. Yu, R. C. Reiner, R. Barber, E. Gaikedu, S. Hay, S. Lim, C. Murray, D. Pigott, B. A. Prakash, B. Adhikari, J. Cui, A. Rodríguez, A. Tabassum, J. Xie, P. Keskinocak, J. Asplund, A. Baxter, B. E. Oruc, N. Serban, S. O. Arik, M. Dusenberry, A. Epshteyn, E. Kanal, L. T. Le, C.-L. Li, T. Pfister, D. Sava, R. Sinha, T. Tsai, N. Yoder, J. Yoon, L. Zhang, S. Abbott, N. I. Bosse, S. Funk, J. Hellewel, S. R. Meakin, J. D. Munday, K. Sherratt, M. Zhou, R. Kalantari, T. K. Yamana, S. Pei, J. Shaman, T. Ayer, M. Adee, J. Chhatwal, O. O. Dalgic, M. A. Ladd, B. P. Linas, P. Mueller, J. Xiao, M. L. Li, D. Bertsimas, O. S. Lami, S. Soni, H. T. Bouardi, Y. Wang, Q. Wang, S. Xie, D. Zeng, A. Green, J. Bien, A. J. Hu, M. Jahja, B. Narasimhan, S. Rajanala, A. Rumack, N. Simon, R. J. Tibshirani, R. Tibshirani, V. Ventura, L. Wasserman, E. B. O’Dea, J. M. Drake, R. Pagano, J. W. Walker, R. B. Slayton, M. Johansson, M. Biggerstaff, and N. G. Reich. Evaluation of individual and ensemble probabilistic forecasts of COVID-19 mortality in the United States. *Proceedings of the National Academy of Sciences*, 119(15):e2113561119, 2022.
- Department of Health and Human Services. COVID-19 guidance for hospital reporting and FAQs for hospitals, hospital laboratory, and acute care facility data reporting, 2023. URL <https://www.hhs.gov/sites/default/files/covid-19-faqs-hospitals-hospital-laboratory-acute-care-facility-data-reporting.pdf>.
- K. Durizzo, K. Harttgen, F. Tediosi, M. Sahu, A. Kuwawenaruwa, P. Salari, and I. Günther. Toward mandatory health insurance in low-income countries? An analysis of claims data in Tanzania. *Health Economics*, 31(10):2187–2207, 2022.
- D. C. Farrow. *Modeling the past, present, and future of influenza*. PhD thesis, Carnegie Mellon University, 2016.
- D. C. Farrow, L. C. Brooks, R. J. Tibshirani, and R. Rosenfeld. Delphi Epidata API, 2015. URL <https://github.com/cmu-delphi/delphi-epidata>.
- J. Geng, X. Chen, J. Shi, H. Bao, Q. Chen, and H. Yu. Assessment of the satisfaction with public health insurance programs by patients with chronic diseases in China: a structural equation modeling approach. *BMC Public Health*, 21(1), 2021.
- R. J. Hyndman and G. Athanasopoulos. *Forecasting: Principles and Practice*. OTexts, third edition, 2021.
- M. Jahja, D. Farrow, R. Rosenfeld, and R. J. Tibshirani. Kalman Filter, sensor fusion, and constrained regression: equivalences and insights. In *Advances in Neural Information Processing Systems*, 2019.
- Y.-S. Jung, Y.-E. Kim, D.-S. Go, R. Munkhzul, J. Jung, and S.-J. Yoon. Associations between private health insurance and medical care utilization for musculoskeletal disorders: using the Korea Health Panel Survey Data for 2014 to 2015. *INQUIRY: The Journal of Health Care Organization, Provision, and Financing*, 57:004695802098146, 2020.
- S. I. Leuba, R. Yaesoubi, M. Antillon, T. Cohen, and C. Zimmer. Tracking and predicting U.S. influenza activity with a real-time surveillance network. *PLOS Computational Biology*, 16(11):1–14, 11 2020.
- S. Li and Y. Yang. An empirical study on the influence of the basic medical insurance for urban and rural residents on Family Financial Asset Allocation. *Frontiers in Public Health*, 9, 2021.

- D. J. McDonald, J. Bien, A. Green, A. J. Hu, N. DeFries, S. Hyun, N. L. Oliveira, J. Sharpnack, J. Tang, R. Tibshirani, V. Ventura, L. Wasserman, and R. J. Tibshirani. Can auxiliary indicators improve COVID-19 forecasting and hotspot prediction? *Proceedings of the National Academy of Sciences*, 118(51):e2111453118, 2021.
- T. Mori, J. Komiyama, T. Fujii, M. Sanuki, K. Kume, G. Kato, Y. Mori, H. Ueshima, H. Matsui, N. Tamiya, and T. Sugiyama. Medical expenditures for fragility hip fracture in Japan: a study using the Nationwide Health Insurance Claims Database. *Archives of Osteoporosis*, 17(1), 2022.
- T. Nakayama, Y. Imanaka, Y. Okuno, G. Kato, T. Kuroda, R. Goto, S. Tanaka, H. Tamura, S. Fukuhara, S. Fukuma, M. Muto, M. Yanagita, and Y. Yamamoto. Analysis of the evidence-practice gap to facilitate proper medical care for the elderly: investigation, using databases, of utilization measures for National Database of Health Insurance Claims and Specific Health Checkups of Japan (NDB). *Environmental Health and Preventive Medicine*, 22(1), 2017.
- T. P. Nickels. RE: CMS-3401-IFC, 2020. URL <https://www.aha.org/system/files/media/file/2020/11/aha-comment-cms-aug-25-interim-final-rule-on-covid-19-data-reporting-letter-11-2-20.pdf>.
- R. Panczak, V. von Wyl, O. Reich, X. Luta, M. Maessen, A. E. Stuck, C. Berlin, K. Schmidlin, D. C. Goodman, M. Egger, K. Clough-Gorr, and M. Zwahlen. Death at no cost? Persons with no health insurance claims in the last year of life in Switzerland. *BMC Health Services Research*, 18(1), 2018.
- J. M. Radin, N. E. Wineinger, E. J. Topol, and S. R. Steinhubl. Harnessing wearable device data to improve state-level real-time surveillance of influenza-like illness in the USA: a population-based study. *The Lancet Digital Health*, 2(2):e85–e93, 2020.
- E. L. Ray, L. C. Brooks, J. Bien, M. Biggerstaff, N. I. Bosse, J. Bracher, E. Y. Cramer, S. Funk, A. Gerding, M. A. Johansson, A. Rumack, Y. Wang, M. Zorn, R. J. Tibshirani, and N. G. Reich. Comparing trained and untrained probabilistic ensemble forecasts of COVID-19 cases and deaths in the United States. *International Journal of Forecasting*, 39(3):1366–1383, 2023.
- A. Reinhart, L. Brooks, M. Jahja, A. Rumack, J. Tang, S. Agrawal, W. A. Saeed, T. Arnold, A. Basu, J. Bien, Á. A. Cabrera, A. Chin, E. J. Chua, B. Clark, S. Colquhoun, N. DeFries, D. C. Farrow, J. Forlizzi, J. Grabman, S. Gratzl, A. Green, G. Haff, R. Han, K. Harwood, A. J. Hu, R. Hyde, S. Hyun, A. Joshi, J. Kim, A. Kuznetsov, W. L. Motte-Kerr, Y. J. Lee, K. Lee, Z. C. Lipton, M. X. Liu, L. Mackey, K. Mazaitis, D. J. McDonald, P. McGuinness, B. Narasimhan, M. P. O’Brien, N. L. Oliveira, P. Patil, A. Perer, C. A. Politsch, S. Rajanala, D. Rucker, C. Scott, N. H. Shah, V. Shankar, J. Sharpnack, D. Shemetov, N. Simon, B. Y. Smith, V. Srivastava, S. Tan, R. Tibshirani, E. Tuzhilina, A. K. V. Nortwick, V. Ventura, L. Wasserman, B. Weaver, J. C. Weiss, S. Whitman, K. Williams, R. Rosenfeld, and R. J. Tibshirani. An open repository of real-time COVID-19 indicators. *Proceedings of the National Academy of Sciences*, 118(51):e2111452118, 2021.
- R. Rosenfeld and R. J. Tibshirani. Epidemic tracking and forecasting: Lessons learned from a tumultuous year. *Proceeding of the National Academy of Sciences*, 118(51):e2111456118, 2021.
- M. Sakai, S. Ohtera, T. Iwao, Y. Neff, G. Kato, Y. Takahashi, and T. Nakayama. Validation of claims data to identify death among aged persons utilizing enrollment data from Health Insurance Unions. *Environmental Health and Preventive Medicine*, 24(1), 2019.
- M. Santillana, A. T. Nguyen, T. Louie, A. Zink, J. Gray, I. Sung, and J. S. Brownstein. Cloud-based electronic health records for real-time, region-specific influenza surveillance. *Scientific Reports*, 6(1):1–8, 2016.
- M. S. Smolinski, A. W. Crawley, K. Baltrusaitis, R. Chunara, J. M. Olsen, O. Wójcik, M. Santillana, A. Nguyen, and J. S. Brownstein. Flu Near You: crowdsourced symptom reporting spanning 2 influenza seasons. *American Journal of Public Health*, 105(10):2124–2130, 2015.
- S. O. Song, E. Han, K. J. Son, B.-S. Cha, and B.-W. Lee. Age at mortality in patients with type 2 diabetes who underwent kidney transplantation: an analysis of data from the Korean National Health Insurance and Statistical Information Service, 2006 to 2018. *Journal of Clinical Medicine*, 12(9), 2023.
- C. Viboud, V. Charu, D. Olson, S. Ballesteros, J. Gog, F. Khan, B. Grenfell, and L. Simonsen. Demonstrating the use of high-volume electronic medical claims data to monitor local and regional influenza activity in the US. *PLOS ONE*, 9(7):1–12, 07 2014.
- C.-Y. Yang, R.-J. Chen, W.-L. Chou, Y.-J. Lee, and Y.-S. Lo. An integrated influenza surveillance framework based on national influenza-like illness incidence and multiple hospital electronic medical records for early prediction of influenza epidemics: design and evaluation. *Journal of Medical Internet Research*, 21(2):e12341, 2019.

- Q. Yao, H. Li, and C. Liu. Use of social health insurance for hospital care by internal migrants in China—Evidence from the 2018 China migrants dynamic survey. *Frontiers in Public Health*, 10, 2022.
- L. Zheng and L. Peng. Effect of major illness insurance on vulnerability to poverty: evidence from China. *Frontiers in Public Health*, 9, 2021.