

Challenges in Real-Time Estimation of Changing Epidemic Severity Rates

Jeremy Goldwasser* Addison Hu* Alyssa Bilinski† Daniel McDonald‡
Ryan Tibshirani*

November 18, 2024

Abstract

Severity rates like case-fatality rate and infection-fatality rate are popular metrics in public health. To guide decision-making in response to changes like new variants or vaccines, it is imperative to understand how these rates shift in real time. We demonstrate that standard ratio estimators for time-varying severity rates may exhibit high statistical bias, failing to detect increases in fatality risk or falsely signalling nonexistent surges. We supplement our theoretical analyses with experimental results on real and simulated COVID-19 data. Finally, we highlight strategies to mitigate this bias, drawing connections with effective reproduction number (R_t) estimation.¹

1 Introduction

Several public health metrics express the probability that a second, more serious outcome will follow a primary event (Garske et al., 2009; Russell et al., 2020b; Challen et al., 2021; Overton et al., 2022; Roth et al., 2021; Xie et al., 2024; Bellan et al., 2020; COVID-19 Forecasting Team, 2022; Luo et al., 2021). For example, the case-fatality rate (CFR) and infection-fatality rate (IFR) are commonly used to assess the intensity of an epidemic. Other examples of such “severity rates” include the hospitalization-fatality rate and case-hospitalization rate.

In an ideal setting, severity rates can be obtained directly from comprehensive line-list or claims data of individual patient outcomes (Roth et al., 2021; Xie et al., 2024; Bellan et al., 2020; Challen et al., 2021). However, in fast-moving epidemics like COVID-19, large-scale tracking is infeasible, especially in real-time (Overton et al., 2022). Instead, rates are routinely estimated from aggregate count data. While many works assume they are constant over time (Reich et al., 2012; Ghani et al., 2005; Jewell et al., 2007; Baud et al., 2020), consequential shifts can occur in response to factors such as new therapeutics, vaccines, and variants (McNeil, 2020). Time-varying severity rates are typically estimated with a ratio of the two aggregate data streams. For example, aggregated cases and deaths were widely used to report COVID CFRs, both in academic literature (Wjst and Wendtner, 2023; Horita and Fukumoto, 2022; Luo et al., 2021; Yuan et al., 2020; Liu et al., 2023) and major news publications like the Atlantic (Madrigal and Moser, 2020) and Wall Street Journal (Kamp and Krouse, 2020). In fact, ratio estimators are so common that IFR, for example, is often referred to as the infection-fatality *ratio* (Luo et al., 2021; COVID-19 Forecasting Team, 2022).

In this work, we demonstrate that these ratio estimators may target the wrong severity rates [DJM: “wrong severity rates” doesn’t sound quite right. It makes me think, for some reason of swapping one for another. Maybe “may target the desired severity rate, but with a bias: the estimates will necessarily be too high or low.”]. This statistical bias arises as a consequence of changing severity rates — precisely when time-varying estimates should be most useful. Bias is influenced by changes in the primary incidence curve (e.g. cases, in CFR), as well as long time delays between events. During COVID-19, we show ratio estimators

*Department of Statistics, University of California, Berkeley

†Brown University School of Public Health

‡Department of Statistics, University of British Columbia

¹Code is available at <https://github.com/jeremy-goldwasser/Severity-Bias>.

would have failed to quickly identify the rise in hospitalization-fatality rate (HFR) during the onset of the Delta wave. After the initial Omicron surge, the severity ratio estimates spiked even though the true HFRs fell. We study the sources of this bias, and suggest alternative methodology which overcomes it.

2 Materials and Methods

[DJM: Some journals don’t mind (and even encourage) empty sections, but I prefer to have some text between section headings. Just a sentence or two is fine.]

2.1 Experimental setup

[DJM: I wonder if it wouldn’t be better to put this section (2.1) as the last section before you move to Results rather than the first?]

[DJM: References. I personally try to use the `cleveref` package so that I can write something like “In `\cref{sec:blah}`...” which gets typeset as “In Section 3...”. You can use `\autoref{sec:blah}` (part of `hyperref` that you’re already using) to get similar results, though it is a bit more fragile. Feel free to do neither. However, 2 things that I think are important: (1) you should do `Section~\ref{sec:blah}` with the tilde to make sure that references don’t get split across lines; (2) equations should be written as either `Equation~\eqref{eq:3}` or `Eq.~\eqref{eq:3}` rather than just `\eqref{eq:3}`. Some journals require one or the other, and I don’t really have a preference, as long as it’s consistent.]

2.1.1 HFR estimation

Our experiments focus on the Hospitalization-Fatality Rate (HFR) during COVID-19. Hospitalization reporting was much more complete than case reporting throughout the pandemic. Hospitals were mandated to report new daily admissions to the Department of Health and Human Services (HHS) or face penalties (Department of Health and Human Services, 2023). The time-to-death delay distribution is indeed supported on integers starting at $k = 0$, since hospitalizations are aligned by admission date.

To estimate real-time HFRs, we use daily hospitalizations and deaths as available in the `epidatr` API, developed by the Delphi Group. Like HHS for hospitalizations (Department of Health and Human Services, 2023), John Hopkins University (JHU) provided the definitive resource for real-time death counts (Center for Systems Science and Engineering (CSSE) at Johns Hopkins University, 2023). These counts reflect the times at which deaths were reported to health authorities, not necessarily when they actually happened. Therefore raw JHU death counts are highly volatile due to reporting idiosyncrasies like day-of-week effects and data dumps. As a result, we used a 7-day trailing average of counts.

Because the daily aggregates from JHU and HHS were updated over the course of the pandemic, we use both the counts available in real time as well as their finalized versions. Occasionally, the most recent date with available counts lagged several days behind the present. To account for this, we estimated HFRs each week, using real-time data that was available two days after each date. In the rare chance that requested counts were still unavailable, we imputed the value with the most recently observed data. Finally, we smoothed counts over the previous 7 days.

The two ratio estimators (2) and (5) require choices of lag and delay distribution. [DJM: I try to avoid referring to future equations. So this is another reason to consider moving this subsection to the end. In fact, I might put this paragraph right at the beginning of Section 3 (before 3.1). It’s a good summary of what comes next there.] Appendix C evaluates the robustness of findings against different hyperparameter values. The experiments in Section 3.1 use a lag of 20 days, which maximizes the cross-correlation between hospitalizations and deaths over all time (Madrigal and Moser, 2020). We let the delay distribution be a discrete gamma, a common choice. We set its mean to this oracle lag, as lags are often chosen to be the mean of the delay distribution. This mean of 20 matches nicely with a UK study that finds a median hospitalization-to-death time of 11 days (Ward and Johnsen, 2021), and a CDC report that 63% of COVID deaths are reported within 10 days (Centers for Disease Control and Prevention, National Center for Health Statistics, 2023). We set the standard deviation to 18, because the delay distributions fit by the UK study had standard deviations that were roughly 90% of their means.

2.1.2 Validation data

While the true HFRs are unknown, there are sound ways to approximate them. One such approach is to use line-list HFRs from the National Hospital Care Survey ([National Center for Health Statistics \(NHCS\), 2023](#)). The NHCS recorded weekly HFRs from inpatient deaths in a representative subset of 601 hospitals across the US.

HFRs from aggregate hospitalization and death counts are significantly higher than those from NHCS because not all deaths occur in hospitals. A CDC analysis reported the percentage of inpatient deaths every month from 2020 through 2022; roughly 60% of COVID deaths occurred in hospitals in 2022, down from nearly 70% in 2021 and 2022. To account for non-inpatient deaths, we divided the NHCS curve by these percentages. Finally, we smoothed the resulting HFRs with a spline. To do so, we used the `smooth.spline` function in R, which chooses the smoothness hyperparameter with generalized cross validation.

We considered two other sources for ground truth HFRs, discussed in [Appendix B.2](#). Unlike NHCS, these HFRs are obtained from aggregate counts, not line-list data. Fortunately, they are fairly consistent with the rescaled NHCS data. Of course, the NHCS curve is merely an approximation for the ground truth; its values, especially after mid-2022, may be incorrect. Nevertheless, it is a useful benchmark to judge the fidelity of our HFR estimates.

2.2 Severity rate estimators

The time-varying severity rate is defined as

$$p_t = \mathbb{P}(\text{secondary event will occur} \mid \text{primary event at time } t). \quad (1)$$

Let $\{X_t\}$, $\{Y_t\}$ denote the time series of interest. In the case of CFR, for example, X_t and Y_t are the total number of new cases and deaths, respectively, at day t .

The canonical estimator for time-varying severity rates is a ratio between X_t and Y_t events, offset by a lag ℓ . This lagged approach is formally introduced in [Thomas and Marks \(2021\)](#), but has also been used in prior works (e.g., [Wjst and Wendtner, 2023](#); [Horita and Fukumoto, 2022](#); [Luo et al., 2021](#); [Yuan et al., 2020](#); [Liu et al., 2023](#); [Madrigal and Moser, 2020](#); [Kamp and Krouse, 2020](#)). The real-time estimator only uses data until the present timestep t :

$$\hat{p}_t^\ell = \frac{Y_t}{X_{t-L}}. \quad (2)$$

Alternative methods use the estimated *delay distribution* that relates the two time series. The delay distribution is defined as [DJM: I found the equation below a bit confusing. I edited, but feel free to change back if you disagree. I think the k means “at time $t+k$ ” not “after k days”. The second case doesn’t seem to rule out that it happens on day $k+j$, $j \geq 0$. This is a (conditional) pmf not a CDF, right? But maybe I’m missing something. The same happens in Equation (3), but I haven’t changed it there. Now that I’ve read through Section 3, is there a cumulative/incidence issue. I think you’re using incidence here and cumulative with γ below. Is that right?]

$$\pi_k^{(t)} := \mathbb{P}(\text{secondary event at } t+k \mid \text{primary event occurs at } t \text{ and secondary event occurs eventually}).$$

Several tools exist to estimate delay distributions from aggregate or line-list data ([Charniga et al., 2024](#)). For ease of analysis, we consider discrete delay distributions, though continuous-time approaches are possible. Similarly, we truncate the delay distribution at d days, in essence assuming all secondary events occur within this period. Finally, we assume delay distributions are constant over time, and suppress the dependence on t in the notation.

The expected number of secondary events at any given day can be expressed in terms of historical primary

incidence, severity rates, and the delay distribution (Qu et al., 2022; Nishiura et al., 2009),²

$$\begin{aligned}
E[Y_t] &= \sum_{k=0}^d X_{t-k} \mathbb{P}(\text{secondary at } t \mid \text{primary at } t-k) \\
&= \sum_{k=0}^d X_{t-k} \mathbb{P}(\text{secondary after } k \mid \text{secondary occurs, primary at } t-k) \\
&\quad \times \mathbb{P}(\text{secondary occurs} \mid \text{primary at } t-k) \\
&= \sum_{k=0}^d X_{t-k} \pi_k p_{t-k}.
\end{aligned} \tag{3}$$

This is a convolution of the delay distribution against the product of primary incidence and the severity rate. If the severity rate is a constant p , (3) simplifies to $E[Y_t] = p \sum_{k=0}^d X_{t-k} \pi_k$. Nishiura et al. (2009) rearranged this expression to estimate this time-stationary rate using a plug-in estimate of the delay distribution and smoothing over the entire history,

$$\hat{p}_t = \frac{\sum_{s=t_0}^t Y_s}{\sum_{s=t_0}^t \sum_{k=0}^d X_{s-k} \gamma_k}. \tag{4}$$

This estimator is widely used in practice (Garske et al., 2009; Russell et al., 2020b,a). Assuming the true rate is indeed stationary and the delay distribution is correctly specified, it is unbiased. Overton et al. (2022) adapted (4) for the time-varying setting, using daily rather than cumulative counts:

$$\hat{p}_t^\gamma = \frac{Y_t}{\sum_{k=0}^d X_{t-k} \gamma_k}. \tag{5}$$

[DJM: What is γ in the equation above? Is that p ? It needs to be defined here.] The convolutional ratio (5) can be understood as a generalization of (2). It reduces to the same ratio that γ is a point mass distribution where all secondary events occur after ℓ days. Otherwise, it may relate the two time series more accurately by means of a smooth delay distribution, since the true distribution is unlikely to be a point mass.

Gupte et al. (2024) used the time-varying convolutional ratio to analyze changing Covid-19 CFRs in the UK. They implemented the convolutional ratios (4) and (5) in the R package `cfr`. In general, however, the lagged ratio is the more commonly used time-varying estimator.

To stabilize estimates, smoothed counts are often used in practice (Wjst and Wendtner, 2023; Luo et al., 2021; Liu et al., 2023). For the sake of simplicity of presentation, we generally focus on the versions described above. However, we formalize the smoothed versions in Equations (9) and (10), and analyze them experimentally.

Severity rates bear natural connections with reproduction numbers. Both the true severity rate as defined in Equation (1) and the case reproduction number R_t are defined as the average number of secondary events produced by a single primary event at t . In contrast to severity rates, reproduction numbers have infections as the primary and secondary events, where a single infection can generate more than one secondary event. Comparable to the delay distribution π is the renewal equation g , measuring the time between primary and secondary infections. [DJM: My understanding is that the renewal equation is actually the whole thing: $y_t = R_t \sum_{k=1}^{\infty} g_k y_{t-k}$. Maybe you mean “generation interval distribution g ”?]

Severity rates and reproduction numbers are also estimated similarly. Because primary events at t produce secondary events after t , their effect is not observed in real time. Therefore, standard real-time estimates for R_t and severity rates (Eq. (2), (5)) analyze the number of *secondary* events at t produced by relevant primary events. For reproduction numbers, this is formally defined as instantaneous R_t , the average number of secondary infections at time t produced by a single primary infection in the past. Indeed, one of the most popular frameworks to estimate R_t is almost identical to the convolutional ratio (Fraser, 2007; Wallinga and Lipsitch, 2007; Cori et al., 2013; Liu et al., 2024). Its only difference is that it uses the same aggregate time

²Throughout this work, we assume primary incidence is known, and condition on $X_{s \leq t}$ implicitly. We also assume the delay distribution π is the same over all time: $\pi_k^{(t)} = \pi_k$ for all k and t .

series, for infections, in both the numerator and denominator:

$$\hat{R}_t = \frac{I_t}{\sum_{k=1}^d I_{t-k} g_k}. \quad (6)$$

2.3 Well-specified analysis

In this section, we explore the bias of the convolutional ratio (5) when the true delay distribution is known.

Theorem 1. *Assume the true delay distribution is a known constant π over all time with maximum length d [DJM: What does maximum length mean here?]. The bias of the convolutional ratio \hat{p}_t^π is*

$$\text{Bias}(\hat{p}_t^\pi) = \sum_{k=0}^d \left[(p_{t-k} - p_t) \frac{X_{t-k} \pi_k}{\sum_{j=0}^d X_{t-j} \pi_j} \right].$$

[DJM: I rewrote this expression just a bit. Since you talk about $(p_{t-k} - p_t)$ first, it may make sense to put it first. Then I added the bracket just so that the sum is a bit more clear. Feel free to revert.] Appendix A contains the short proof.

The degree of bias in Theorem 1 depends on three factors.

1. **Changes in severity rate.** The central component of this bias expression is the difference $(p_{t-k} - p_t)$. When severity rates are constant over the d preceding days, this estimator is unbiased (because this difference is zero). This is in line with the unbiasedness of the estimator using cumulative counts assuming a globally stationary rate (Nishiura et al., 2009). But when severity rates change before t , these difference terms will be nonzero, in which case the estimator will be biased.³ Figures 1a and 13a illustrate this scenario: the estimated severity rates are most inaccurate at periods where the true rate is changing quickly.
- To make matters worse, the bias is in the opposite direction of the trend we want to detect. For example, suppose the severity rate is monotonically falling, with $p_t < p_{t-1} < \dots < p_{t-d}$. As a result, the bias is positive, meaning the ratio estimates do not decline with the true rate. In fact, the estimated severity may even rise, not fall. Conversely, when true severity rates are rising, the ratio estimates will be too low.
2. **The delay distribution.** How much the changing severity rates impact the bias depends on the shape of the delay distribution π . In general, the bias is greatest when the delay distribution is long-tailed enough to upweight significant differences in severity rate. While this distinction may appear subtle, Section 3 highlights its surprisingly large effects. The simple example in Figures 1b and 13a shows significant differences in bias between shorter and longer delay distributions.
3. **The primary incidence curve.** Changing primary incidence will also affect the bias, presuming the severity rate changes roughly monotonically in the recent past. Intuitively, this up- or down-weights the terms $(p_{t-k} - p_t) X_{t-k} \pi_k$ for dates further from the present, which are likely to contribute the most bias. In general, falling primary incidences will amplify the bias, whereas rising events will minimize it; see Appendix D.2 for more detail. Figures 1c and 13b visualize this trend on the convolutional ratio.

As noted previously, the convolutional ratio is equivalent to the lagged ratio if γ is a point mass distribution at ℓ . In this oracle setting, all secondary events occur after exactly ℓ days, a highly unrealistic situation. Nevertheless, if this is the case, then

$$\text{Bias}(\hat{p}_t^\ell) = \text{Bias}(\hat{p}_t^\gamma) = p_{t-\ell} - p_t.$$

This toy setting and others are discussed in Appendix D.2.

³It is possible that this estimator could still be unbiased in the unlikely event that individual components in the summation over k exactly cancel each other.

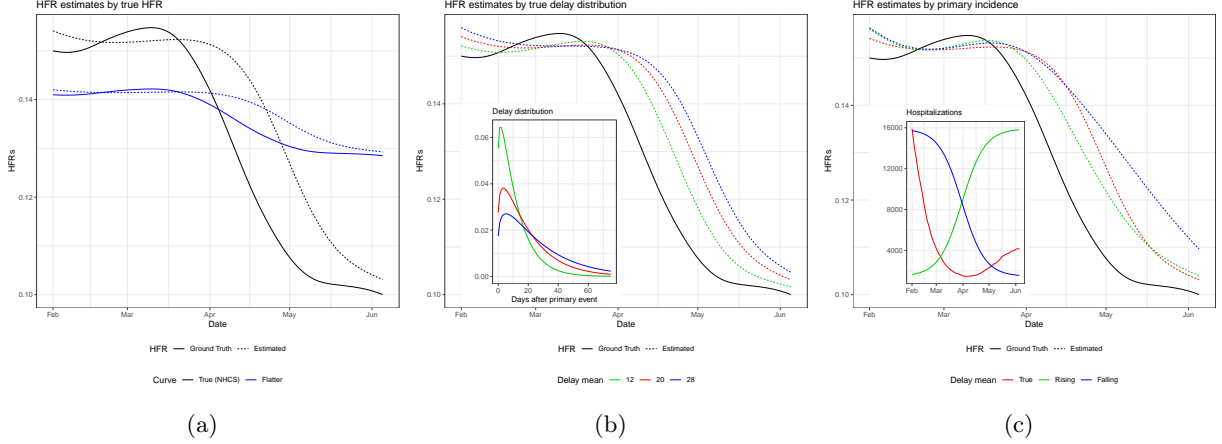


Figure 1: Simple examples of severity rate bias to illustrate the three factors. Deaths were computed noiselessly from (3), with NCHS HFRs and HHS hospitalizations during early 2022. Panels 1a and 1c use the delay distribution estimated with JHU deaths (Section 2.1).

2.4 Misspecified analysis

The above section considered the bias of the convolutional ratio where the true delay distribution π is known. We now consider the more general case, in which it is instead replaced with a plug-in estimate γ . Note the bias of the lagged estimator is a special case, where the plug-in distribution is a point mass at lag time ℓ .

Theorem 2. *Assume that the true delay distribution π is constant over time and that its maximal length d exceeds that of the plug-in distribution γ . Define $A_t^\gamma := \sum_{j=0}^d X_{t-j}\pi_j / \sum_{j=0}^d X_{t-j}\gamma_j$, which compares how the delay distributions convolve against the most recent primary incidence levels. The misspecified bias is*

$$\text{Bias}(\hat{p}_t^\gamma) = A_t^\gamma \text{Bias}(\hat{p}_t^\pi) + p_t(A_t^\gamma - 1).$$

Theorem 2, proven in Appendix A, provides an additive decomposition of the misspecified ratio’s bias. In both terms, A_t^γ dictates the extent to which the misspecified distribution alters the bias. The first term scales the oracle bias, whereas the other solely expresses misspecification. Which of these two terms will dominate depends on the true severity rate p_t , the oracle bias, and the ratio A_t^γ . If oracle bias is small, for example, then its multiplicative scaling should have relatively little effect, in which case the misspecification term may drive bias. This seems to generally be the case in the simulated experiments in Section 3.2.

[DJM: I think that Figure 2 needs a direct discussion somewhere in these next few paragraphs. Like “Figure 2 illustrates...”] Suppose primary events have stabilized after falling for a long time (see April 2022 in Fig. 2). If the plug-in delay distribution is too light-tailed, then $A_t^\gamma > 1$, since it does not upweight distant dates with high primary counts. Therefore, this distribution inflates the oracle bias multiplicatively and adds positive misspecification bias. If primary events have consistently risen instead, then $A_t^\gamma < 1$, so the oracle bias term would shrink and the misspecification bias would be negative. This explains the negative bias of the light-tailed distribution in mid-January. These relations may be more complicated if primary incidence has changed direction throughout the delay distribution.

For the lagged estimator, Equation (2) becomes

$$\text{Bias}(\hat{p}_t^\ell) = \frac{\sum_{j=0}^d X_{t-j}\pi_j}{X_{t-\ell}} \text{Bias}(\hat{p}_t^\pi) + p_t \left(\frac{\sum_{k=0}^d X_{t-k}\pi_k}{X_{t-\ell}} - 1 \right). \quad (7)$$

In some cases, $\gamma_k = \mathbf{1}\{k = \ell\}$ can be thought as a light-tailed distribution. It assigns all its mass at ℓ — chosen to be around the mean of π — and none in the long tail of π . In the flattened-out period around April 2022, the lagged estimator has similar positive bias as the short-tailed convolutional ratio.

Overall, however, the lagged estimator’s bias is more subtle because it relies exclusively on primary incidence ℓ days ago. When counts rise sharply between $t - \ell$ and t , then $A_t^\ell > 1$ due to its small denominator.

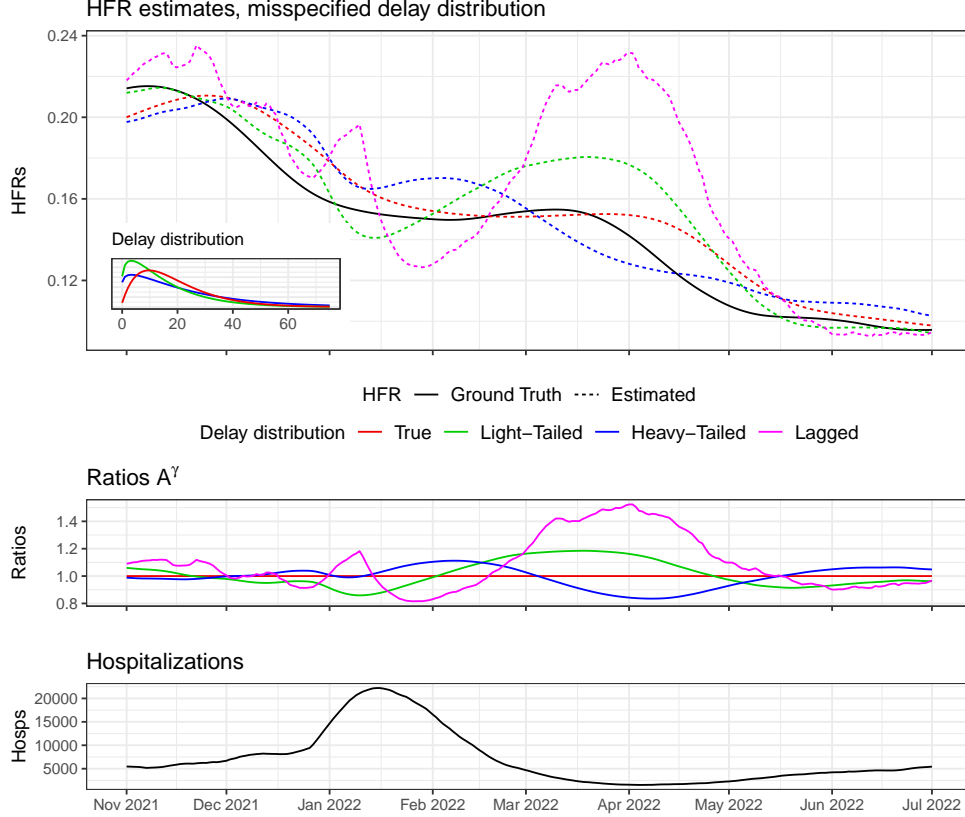


Figure 2: HFR estimates under misspecification. HHS hospitalizations, NCHS HFRs, and noiseless deaths from Eq. (3). Convolutional ratio estimates with true delay distribution (mean 20), misshapen gammas (mean 16 and 24), and point mass at correlation-maximizing lag (16).

This contrasts with the smooth, light-tailed γ described above, which emphasizes recent high counts. Figure 2 highlights this divergence in behavior as hospitalizations peak in mid-January.

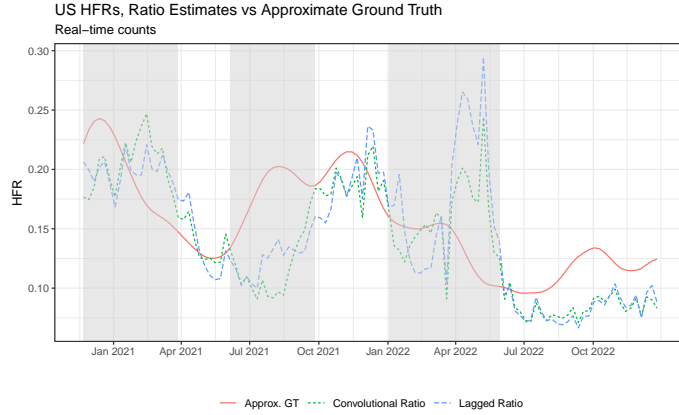
The lagged estimator also has interesting behavior as primary incidence falls from its peak. The denominator of A_t^ℓ is large, since counts neared the peak ℓ days ago. Meanwhile, the numerator is smaller due to its inclusion of lower counts before and after the peak. As a result, $A_t^\ell < 1$, contributing negative bias. Misspecified smooth delay distributions will be less biased under these conditions, since they incorporate the lower counts into the denominator of A_t^γ . This accounts for the lagged ratio’s spurious dip in February 2022.

Note the denominator of A_t^γ , $\sum_{j=0}^d X_{t-j} \gamma_j$, is most extreme when γ is a point-mass distribution at the maximal or minimal values of X_{t-d}, \dots, X_t . This inevitably occurs with the lagged estimator as its convolution sweeps across the primary incidence curve. In contrast, the convolutional ratios do not reach such extremes by nature of their smooth delay distributions. This explains why in Figure 2, the ratio A_t^ℓ fluctuates furthest away from 1 for the lagged estimator. Hence, the lagged ratio can occasionally be understood as providing the worst-case delay distribution in Theorem 2.

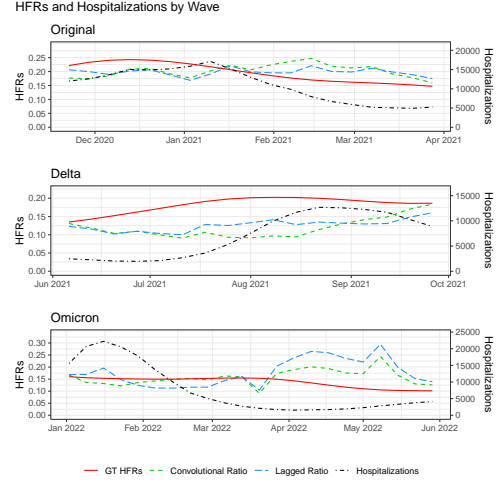
3 Results

3.1 National COVID data

Figure 3 highlights the bias of these ratio estimators [DJM: as described in Theorems 1 and 2 when applied to...]. Both the lagged and convolutional ratios respond very slowly to changes in the HFR. As the HFR declines following the wave in winter 2021, both ratios remain near 0.2 for several months. More troublingly,



(a) Comparing convolutional and lagged ratios against approximate ground truth.



(b) HFRs and hospitalizations in three periods with major bias.

Figure 3: Convolutional ratio estimates are biased regardless of which delay distribution is selected. Real-time counts, Nov. 2020 - Dec 2022. Biased periods of major waves are highlighted.

they are very slow to detect the rising HFR in the early Delta period (summer 2021). If the purpose of these estimators is to inform stakeholders of increased risks in real time, they fail during the Delta surge.

The most significant bias comes in the middle of the Omicron wave in spring 2022. In this period, the HFR remains around 15% until April, then sharply declines to 9% two months later. The lagged ratios first oscillate above and below the true HFRs. Subsequently, both estimates surge as the true HFR nears its nadir, with the lagged ratio nearing 30%. This dramatic upswing signals a serious false alarm. The analysis in Sections 2.3 and 2.4 explain each of these failure cases. We start by analyzing the convolutional ratio with respect to the well-specified bias expression in Theorem 1. While this expression assumes that the true delay distribution is known, we found that different choices of delay distribution generally yield the same bias (Appendix C). This indicates that our estimates may not be far from the oracle ratio.

First, Theorem 1 indicates that the bias moves in the opposite direction of the true severity rate. This occurs during the Delta wave, when the HFRs rise well before the ratio estimates do. Conversely, falling HFRs produce positive bias, as observed in the original and Omicron waves. Second, the enormity of the bias during Omicron can partially be attributed to the precipitous decline in hospitalizations, as falling primary incidence has been shown to exacerbate the bias. Average daily hospitalizations declined from over 20,000 in mid-January to only 1,500 by April 1. Finally, the delay distribution is relatively long with JHU deaths due to its alignment by report date. This is shown to have a substantial impact on the bias, as analyzed in Appendix B.1. Third, the misspecification analysis explains central discrepancies between the convolutional and lagged ratios. The lagged ratio’s erratic behavior during Omicron is almost identical to the simulated results in Figure 2. This bias is likely due to changes in hospitalization counts, which affect A_t^ℓ in Theorem 2. Finally, the lagged ratio is briefly less negatively biased in August 2021, during the Delta wave. Like the spike in January 2022, this can be explained by the sharp rise in hospitalizations. The surge causes A_t^ℓ to rise above 1, so the additive term in Eq. 2 is positive. This briefly offsets the negative oracle bias, bringing the lagged ratio closer to the true HFR. Similarly, the lagged bias is less *positive* in February 2021. Here, $A_t^\ell < 1$ due to falling hospitalizations, so the additive term is negative.

We performed several robustness checks to assess the stability of these findings. Appendix C explores the effect of different hyperparameters and locations. By and large, the ratio estimators yield roughly the same bias regardless of these considerations. It also compares HFR estimates using finalized counts, rather than the data available in real time. This exploration finds that the observed biases could not be attributed to real-time reporting issues.

HFRs, Simulated Deaths

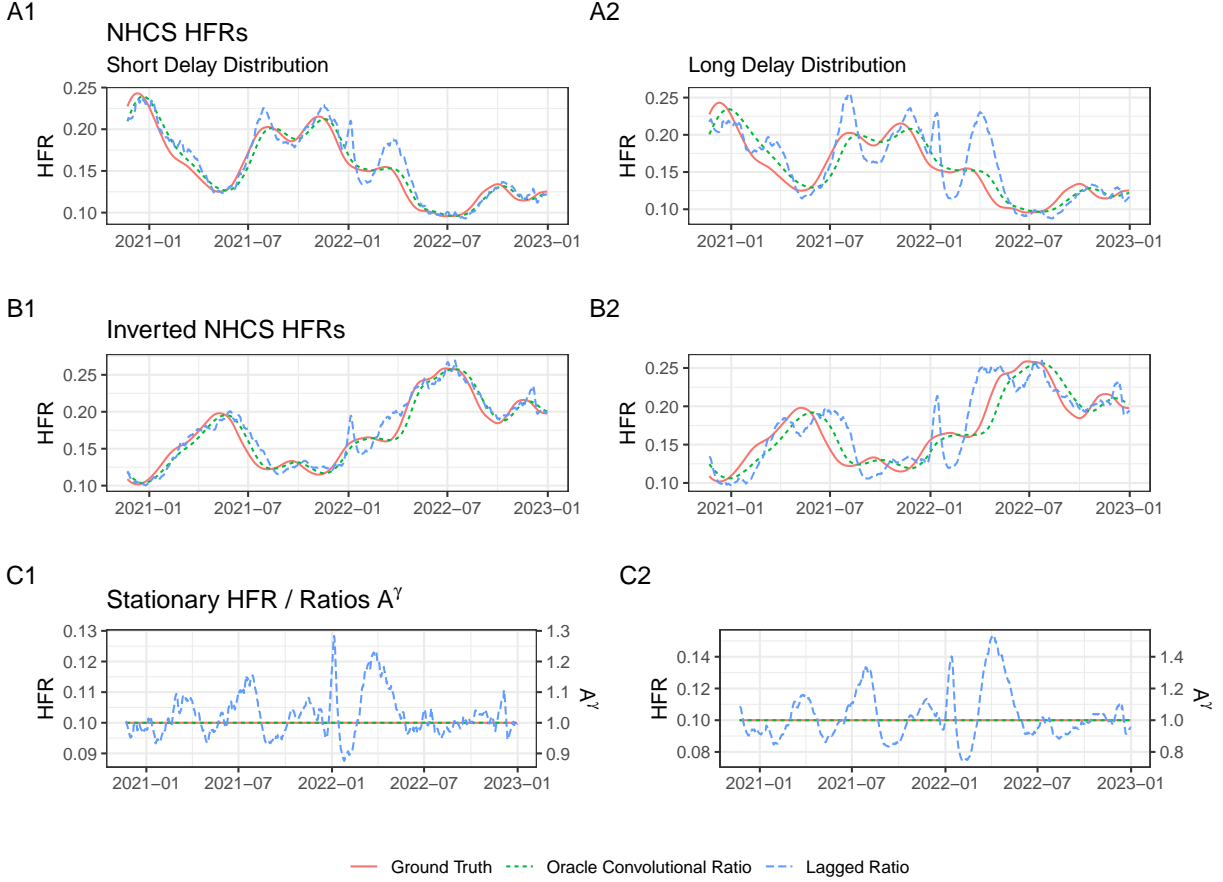


Figure 4: True and Estimated HFRs from Simulated Deaths. First column has short delay distribution, second has long.

3.2 Simulated data

We further evaluated these methods in a variety of simulation settings. Given a series of time-varying HFRs p_t and delay distribution π , deaths are defined without noise from (3):

$$Y_t := \sum_{k=0}^d X_{t-k} \mathbb{P}(\text{die at } t \mid \text{hosp at } t-k) = \sum_{k=0}^d X_{t-k} \pi_k p_{t-k}.$$

Like the experiments in Sections 2.3 and 2.4, we used finalized HHS hospitalization counts. However, the simulations in this section evaluate performance over a two-year period, with a broader range of underlying HFR curves. To supplement the NHCS HFRs, we mimicked the opposite trend by inverting and rescaling them. We also modeled a stationary HFR of 10% over all time. As in Section 2.1, the delay distributions were again gamma with standard deviation 0.9 of their mean. We experimented with means of 12 and 24 to illustrate a short and long delay distribution.

To elucidate the oracle bias in Theorem 1, we let the convolutional ratio use the true delay distribution. For the lagged ratio, ℓ was again chosen to maximize the cross-correlation between hospitalizations and deaths. We also experimented with the mean of the delay distribution, proposed in Feng et al. (2023). Figure 12 in Appendix D shows the results are similarly unstable.

Figure 4 displays the results on the six settings of delay distribution and HFR. In addition to the HFR curves, it also visualizes the ratios A_t^γ , which merely rescale the HFR estimates in the stationary case. Given a

constant rate p , the oracle convolutional ratio is unbiased, so Theorem 2 reduces to $\mathbb{E}[\hat{p}_t^\gamma] = pA_t^\gamma$. Furthermore, $\mathbb{E}[\hat{p}_t^\gamma] = \hat{p}_t^\gamma$, as our setup simulates deaths without noise. Consequently, $A_t^\gamma = \frac{\hat{p}_t^\gamma}{p}$.

Matching expectations, HFRs are significantly more biased given the longer underlying delay distribution. In all three HFR settings, the lagged ratios swing more widely. For example, when the true HFR is a constant 10%, they peak at 12% under the light-tailed distribution, compared to 15% with the heavy tail. The oracle convolutional ratio does not share the lagged estimator’s dramatic oscillations. Rather, it tracks the general shape of true curve, albeit at a delay. For the NHCS HFRs, the average delay was 5 days for the light-tailed distribution, and 12 days with heavy tail; these delays were 6 and 15 days for the inverted curve. (To compute the average delay, we again took the maximal cross-correlation between the two series.)

In general, the oracle convolutional ratio performs much better than the lagged estimator. Our analysis accounts for this wide gap in performance. In the case of a stationary severity rate, for example, Theorem 1 assures the oracle convolutional ratio is unbiased. But if the delay distribution is misspecified, Theorem 2 expresses the resulting bias via the ratio A_t^γ . This quantity does not depend on the severity rate, so the bias moves in similar trajectories across the three HFR settings.

The analysis in Section 2.4 explains these trajectories. During the Delta and Omicron waves, rapid rises in hospitalizations produced high values of A_t^ℓ . This accounts for the spikes in August 2021 and January 2022, across all 3 HFR settings. Positive bias is also expected when hospitalizations have leveled out from a decline. This consistently occurs in spring 2022, with A_t^ℓ reaching 1.2 and 1.5 for the short and long distributions. Lastly, the lagged estimator should have negative bias as primary events fall. We observe this in Delta (September 2021) and Omicron (February 2022).

The misspecified bias (Theorem 2) rescales the oracle bias and adds a misspecification term. Studying Figure 4, we observe the misspecification term tends to dominate when A_t^γ strays away from 1. To understand this, consider periods in which the oracle bias is negative. Here, the oracle and misspecification terms are at odds with each other. $A_t^\gamma > 1$ amplifies the negative bias while adding positive misspecification bias, whereas $A_t^\gamma < 1$ has the opposite effect.

Invariably, the lagged ratio’s bias moves in the direction of the misspecification term $p_t(A_t^\gamma - 1)$. Under the true NHCS HFRs, for example, the lagged estimates spike with A_t^ℓ in August 2021. In the inverted setting, the lagged bias tracks the down-up-down motion of A_t^ℓ during the first five months of 2021. That the misspecification term wins out in these conflicting settings indicates it accounts for a disproportionate amount of the bias.

For a further example, consider the bias at April 2022 in panel A2. The true HFR is 14%, with the convolutional ratio nearby at 15%. Meanwhile, the lagged ratio peaks at 22.5%, driven upwards by an A_t^ℓ of 1.5. Decomposing the lagged bias of 8.5% with Theorem 2, the oracle term $A_t^\ell \text{Bias}(\hat{p}_t^\pi)$ equals only 1.5%; meanwhile, the misspecification term $p_t(A_t^\ell - 1) = 7\%$, accounting for the majority of the bias.

4 Discussion

Our analyses illustrate that practitioners should take caution when using time-varying severity ratio estimators. They exhibit considerable bias when severity rates change, particularly the popular lagged ratio estimator. A major purpose of these estimators is to inform stakeholders of changing risks in real time; this bias indicates they may fail to do so in a reliable manner.

Analyzing the lagged ratio enables us to make real-time heuristics about its performance in practice. Theorem 2 decomposes its bias into oracle and misspecification terms, the latter of which has been shown to dominate. Based solely on the primary incidence curve, we can expect the lagged ratio to make the following errors:

1. Unreasonably high severity estimates when primary incidence is rising quickly;
2. Rapid declines when primary incidence is falling quickly;
3. Unexpected surges when primary incidence has leveled out after falling.

Practitioners can adjust their reactions accordingly when these bias patterns occur in real time. For example, if the lagged HFR spikes shortly after hospitalizations reach a stable low, a savvy epidemiologist can temper her alarm with the knowledge it may well be spurious.

While the lagged estimator is ubiquitous in practice, our analysis of its drawbacks suggests other aggregate estimators should be favored. Figure 4 showed the convolutional ratio has the capability to be much more accurate. Even with a rough estimate of the delay distribution, it generally displayed improved stability and performance (Figures 2, 3, and 4). While Equation (4) from Nishiura et al. (2009) is widely used to compute stationary or average HFRs, we have not come across any applications for the time-varying case.

Qu et al. (2022) propose an alternative approach which differs sharply from the ratios discussed in this paper. Their method estimates all historical severity rates at once, using the relation in (3) to fit a Fused Lasso model. This estimator is inherently forward-looking, where rates at t are exclusively used to produce secondary events after t . In spite of these advantages, it may suffer from other sources of bias. It is inclined to estimate smoothly-changing severity rates as piecewise constant, and may yield unstable real-time estimates due to scarce data at the tail. Thorough investigation of its performance is a promising object for future study. Overton et al. (2022) also proposed a forward-looking method, this one a ratio between relevant primary and secondary events. However, this method is not applicable in real time, as it uses secondary events after t to compute the severity rate. Nevertheless, it is a useful tool for retrospective estimation.

Another retrospective tool is aggregate COVID deaths from NCHS, a resource that was not available in real time (Appendix B.1). Unlike JHU, whose aggregates align deaths by report date, NCHS counts deaths on the day the actually occurred. As a result, the mean of its delay distribution is considerably lower, so it produces more accurate ratio estimates (Figures 4 and 5). Analogously, bias is a more serious issue with earlier primary events. For example, case- or infection-fatality ratios may be more biased than hospitalization-fatality ratios.

Severity rates may be biased in ways beyond the statistical bias our work focuses on. Section 2.1 mentioned, for example, the fact that estimating HFR from aggregates fails to address the large proportion of deaths that occur outside the hospital; Lipsitch et al. (2015) refers to this as “survivorship bias.” A central challenge for CFR estimation is under-reporting: Not all events are reported, reporting rates change across time, and severe cases are more likely to be reported than mild cases. Reich et al. (2012) proposes an estimator for a time-invariant *relative* CFR - the ratio of CFRs between groups - that learns these latent reporting rates via the EM algorithm (Dempster et al., 1977). Angelopoulos et al. (2020) applied this in the context of COVID-19, analyzing how the chosen delay distribution affects its results. They also identify other sources of bias, like differences in case definition and testing eligibility.

As discussed in Section 2.2, severity rates may be understood in connection with reproduction numbers. This connection extends to their bias as well. For example, we demonstrated that the convolutional ratio (5) is unbiased if the severity rate and delay distribution in the d days before t are stationary. In a similar vein, Fraser (2007) notes that instantaneous R_t is equal to case R_t if conditions remain unchanged. Future work along the lines of Eales and Riley (2023) could apply our analytical framework to R_t bias, examining the fidelity of instantaneous R_t as a proxy for case R_t .

References

- Adjei, Stacey and Hong, Kai and Molinari, Noelle-Angelique M and Bull-Otterson, Lara and Ajani, Umed A and Gundlapalli, Adi V and Harris, Aaron M and Hsu, Joy and Kadri, Sameer S and Starnes, Jon and Yeoman, Kristin and Boehmer, Tegan K (2022). Mortality risk among patients hospitalized primarily for covid-19 during the omicron and delta variant pandemic periods - united states, april 2020-june 2022. *MMWR Morb Mortal Wkly Rep*, 71(37):1182–1189.
- Angelopoulos, A. N., Pathak, R., Varma, R., and Jordan, M. I. (2020). On Identifying and Mitigating Bias in the Estimation of the COVID-19 Case Fatality Rate. *Harvard Data Science Review*, (Special Issue 1). <https://hdrs.mitpress.mit.edu/pub/y9vc2u36>.
- Baud, D., Qi, X., Nielsen-Saines, K., Musso, D., Pomar, L., and Favre, G. (2020). Real estimates of mortality following COVID-19 infection. *Lancet Infect Dis*, 20(7):773. Epub 2020 Mar 12.
- Bellan, M., Patti, G., Hayden, E., et al. (2020). Fatality rate and predictors of mortality in an italian cohort of hospitalized covid-19 patients. *Sci Rep*, 10:20731.

- Center for Systems Science and Engineering (CSSE) at Johns Hopkins University (2023). Covid-19 data repository. GitHub repository.
- Centers for Disease Control and Prevention, National Center for Health Statistics (2023). Deaths by select demographic and geographic characteristics. Archived September 27, 2023.
- Challen, R., Brooks-Pollock, E., Read, J. M., Dyson, L., Tsaneva-Atanasova, K., and Danon, L. (2021). Risk of mortality in patients infected with SARS-CoV-2 variant of concern 202012/1: matched cohort study. *BMJ*, 372:n579.
- Charniga, K., Park, S. W., Akhmetzhanov, A. R., Cori, A., Dushoff, J., Funk, S., Gostic, K. M., Linton, N. M., Lison, A., Overton, C. E., Pulliam, J. R. C., Ward, T., Cauchemez, S., and Abbott, S. (2024). Best practices for estimating and reporting epidemiological delay distributions of infectious diseases using public health surveillance and healthcare data.
- Cori, A., Ferguson, N. M., Fraser, C., and Cauchemez, S. (2013). A new framework and software to estimate time-varying reproduction numbers during epidemics. *American Journal of Epidemiology*, 178(9):1505–1512.
- COVID-19 Forecasting Team (2022). Variation in the COVID-19 infection–fatality ratio by age, time, and geography during the pre-vaccine era: a systematic analysis. *The Lancet*, 399(10334):1469–1488.
- Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1):1–38.
- Department of Health and Human Services (2023). Covid-19 guidance for hospital reporting and faqs for hospitals, hospital laboratory, and acute care facility data reporting.
- Eales, O. and Riley, S. (2023). Differences between the true reproduction number and the apparent reproduction number of an epidemic time series.
- Feng, J., Luo, H., Wu, Y., Zhou, Q., and Qi, R. (2023). A new method for accurate calculation of case fatality rates during a pandemic: Mathematical deduction based on population-level big data. *Infectious Medicine*, 2(2):96–104.
- Fraser, C. (2007). Estimating Individual and Household Reproduction Numbers in an Emerging Epidemic. *PLOS ONE*, 2(8):1–12.
- Garske, T., Legrand, J., Donnelly, C. A., Ward, H., Cauchemez, S., Fraser, C., Ferguson, N. M., and Ghani, A. C. (2009). Assessing the severity of the novel influenza a/h1n1 pandemic. *BMJ*, 339.
- Ghani, A. C., Donnelly, C. A., Cox, D. R., Griffin, J. T., Fraser, C., Lam, T. H., Ho, L. M., Chan, W. S., Anderson, R. M., Hedley, A. J., and Leung, G. M. (2005). Methods for Estimating the Case Fatality Ratio for a Novel, Emerging Infectious Disease. *American Journal of Epidemiology*, 162(5):479–486.
- Gupte, P., Kucharski, A., Russell, T., Lambert, J., Gruson, H., Taylor, T., Azam, J., Degoot, A., and Funk, S. (2024). cfr: Estimate disease severity and case ascertainment. *Data Collection*. Comprehensive R Archive Network. <https://cran.r-project.org/package=cfr>.
- Horita, N. and Fukumoto, T. (2022). Global case fatality rate from COVID-19 has decreased by 96.8% during 2.5 years of the pandemic. *Journal of Medical Virology*.
- Jewell, N. P., Lei, X., Ghani, A. C., Donnelly, C. A., Leung, G. M., Ho, L.-M., Cowling, B. J., and Hedley, A. J. (2007). Non-parametric estimation of the case fatality ratio with competing risks data: an application to Severe Acute Respiratory Syndrome (SARS). *Stat Med*, 26(9):1982–1998.
- Kamp, J. and Krouse, S. (2020). Case-Fatality Metric Points to Increase in December Deaths. *Wall Street Journal*.
- Lipsitch, M., Donnelly, C. A., Fraser, C., Blake, I. M., Cori, A., Dorigatti, I., Ferguson, N. M., Garske, T., Mills, H. L., Riley, S., Van Kerkhove, M. D., and Hernán, M. A. (2015). Potential biases in estimating absolute and relative case-fatality risks during outbreaks. *PLoS Neglected Tropical Diseases*, 9(7):e0003846.

- Liu, J., Cai, Z., Gustafson, P., and McDonald, D. J. (2024). Time-varying reproduction number estimation with trend filtering. *PLOS Computational Biology*.
- Liu, J., Wei, H., and He, D. (2023). Differences in case-fatality-rate of emerging sars-cov-2 variants. *Public Health in Practice*, 5:100350.
- Luo, G., Zhang, X., Zheng, H., and He, D. (2021). Infection fatality ratio and case fatality ratio of covid-19. *International Journal of Infectious Diseases*, 113:43–46.
- Madrigal, A. C. and Moser, W. (2020). How Many Americans Are About to Die? *The Atlantic*.
- McNeil, D. G. J. (2020). The Pandemic’s Big Mystery: How Deadly Is the Coronavirus? *New York Times*.
- National Center for Health Statistics (NCHS) (2023). In-hospital mortality among hospital confirmed covid-19 encounters by week from selected hospitals. National Hospital Care Survey (NHCS).
- Nishiura, H., Klinkenberg, D., Roberts, M., and Heesterbeek, J. A. P. (2009). Early Epidemiological Assessment of the Virulence of Emerging Infectious Diseases: A Case Study of an Influenza Pandemic. *PLoS One*, 4(8):e6852.
- Overton, C., Webb, L., Datta, U., Fursman, M., Hardstaff, J., Hiironen, I., Paranthaman, K., Riley, H., Sedgwick, J., Verne, J., Willner, S., Pellis, L., and Hall, I. (2022). Novel methods for estimating the instantaneous and overall COVID-19 case fatality risk among care home residents in England. *PLoS Comput Biol*, 18(10):e1010554.
- Qu, Y., Lee, C. Y., and Lam, K. F. (2022). A novel method to monitor covid-19 fatality rate in real-time, a key metric to guide public health policy. *Sci Rep*, 12:18277.
- Reich, N. G., Lessler, J., Cummings, D. A. T., and Brookmeyer, R. (2012). Estimating Absolute and Relative Case Fatality Ratios from Infectious Disease Surveillance Data. *Biometrics*, 68(2):598–606. Published online 2012 Jan 25.
- Roth, G. A., Emmons-Bell, S., Alger, H. M., Bradley, S. M., Das, S. R., de Lemos, J. A., Gakidou, E., Elkind, M. S. V., Hay, S., Hall, J. L., Johnson, C. O., Morrow, D. A., Rodriguez, F., Rutan, C., Shakil, S., Sorensen, R., Stevens, L., Wang, T. Y., Walchok, J., Williams, J., and Murray, C. (2021). Trends in Patient Characteristics and COVID-19 In-Hospital Mortality in the United States During the COVID-19 Pandemic. *JAMA Network Open*, 4(5):e218828–e218828.
- Russell, T. W., Hellewell, J., Abbott, S., Jarvis, C. I., van Zandvoort, K., Ratnayake, R., CMMID nCov working group, Flasche, S., Eggo, R., Edmunds, W. J., and Kucharski, A. J. (2020a). Using a Delay-Adjusted Case Fatality Ratio to Estimate Under-Reporting. *Fondazione Cerm*.
- Russell, T. W., Hellewell, J., Jarvis, C. I., van Zandvoort, K., Abbott, S., Ratnayake, R., working group, C. C., Flasche, S., Eggo, R. M., Edmunds, W. J., and Kucharski, A. J. (2020b). Estimating the infection and case fatality ratio for coronavirus disease (COVID-19) using age-adjusted data from the outbreak on the Diamond Princess cruise ship, February 2020. *Eurosurveillance*, 25(12).
- Thomas, B. S. and Marks, N. A. (2021). Estimating the case fatality ratio for covid-19 using a time-shifted distribution analysis. *Epidemiol Infect*, 149:e197.
- Wallinga, J. and Lipsitch, M. (2007). How generation intervals shape the relationship between growth rates and reproductive numbers. *Proceedings of the Royal Society B: Biological Sciences*, 274(1609):599–604.
- Ward, T. and Johnsen, A. (2021). Understanding an evolving pandemic: An analysis of the clinical time delay distributions of covid-19 in the united kingdom. *PLoS One*, 16(10):e0257978.
- Wjst, M. and Wendtner, C. (2023). High variability of COVID-19 case fatality rate in Germany. *BMC Public Health*, 23:416.

- Xie, Y., Choi, T., and Al-Aly, Z. (2024). Mortality in Patients Hospitalized for COVID-19 vs Influenza in Fall-Winter 2023-2024. *JAMA*, 331(22):1963–1965.
- Yuan, J., Li, M., Lv, G., and Lud, Z. K. (2020). Monitoring transmissibility and mortality of COVID-19 in Europe. *Int J Infect Dis*, 95:311–315.

A Bias Proofs

We first prove Theorem 1, the bias of the oracle convolutional ratio.

Proof.

$$\begin{aligned}
\text{Bias}(\hat{p}_t^\pi) &= E[\hat{p}_t^\pi] - p_t \\
&= \frac{E[Y_t]}{\sum_{k=0}^d X_{t-k}\pi_k} - p_t \\
&= \frac{\sum_{k=0}^d X_{t-k}\pi_k p_{t-k}}{\sum_{k=0}^d X_{t-k}\pi_k} - \frac{p_t \sum_{k=0}^d X_{t-k}\pi_k}{\sum_{k=0}^d X_{t-k}\pi_k} \\
&= \sum_{k=0}^d \frac{X_{t-k}\pi_k}{\sum_{j=0}^d X_{t-j}\pi_j} (p_{t-k} - p_t).
\end{aligned}$$

□

The well-specified bias can be understood as a weighted average of $\{p_{t-k} - p_t\}_{k=0}^d$. The attainable absolute bias ranges between $\min_{k=0,\dots,d} |p_{t-k} - p_t| = 0$, achieved by $k = 0$, and $\max_{k=0,\dots,d} |p_{t-k} - p_t|$. This maximal bias is achieved by setting one of the weights $X_{t-k}\pi_k / (\sum_{j=0}^d X_{t-j}\pi_j)$ to 1 and the rest to zero, either through the delay distribution π or through the primary incidence curve X . Hence, the explanations for delay distribution and primary incidence are aligned: They inflate the bias by upweighting distant timepoints for which the severity rate was different. If severity rates are monotonically changing, for example, then the maximal bias occurs at $k = d$.

Next, we prove the bias expression under misspecified delay distribution (Theorem 2).

Proof.

$$\begin{aligned}
\text{Bias}(\hat{p}_t^\gamma) &= \frac{E[Y_t]}{\sum_{k=0}^d X_{t-k}\gamma_k} - p_t \\
&= \frac{\sum_{k=0}^d X_{t-k}\pi_k p_{t-k}}{\sum_{k=0}^d X_{t-k}\gamma_k} - \frac{\sum_{k=0}^d X_{t-k}\gamma_k p_t}{\sum_{k=0}^d X_{t-k}\gamma_k} \\
&= \sum_{k=0}^d \frac{X_{t-k}}{\sum_{j=0}^d X_{t-j}\gamma_j} (\pi_k p_{t-k} - \gamma_k p_t) \\
&= \sum_{k=0}^d \frac{X_{t-k}}{\sum_{j=0}^d X_{t-j}\gamma_j} (\pi_k p_{t-k} - (\pi_k + (\gamma_k - \pi_k)) p_t) \\
&= \frac{\sum_{j=0}^d X_{t-j}\pi_j}{\sum_{j=0}^d X_{t-j}\gamma_j} \sum_{k=0}^d \frac{X_{t-k}\pi_k}{\sum_{j=0}^d X_{t-j}\pi_j} (p_{t-k} - p_t) - \\
&\quad p_t \sum_{k=0}^d \frac{X_{t-k}}{\sum_{j=0}^d X_{t-j}\gamma_j} (\gamma_k - \pi_k) \\
&= \frac{\sum_{j=0}^d X_{t-j}\pi_j}{\sum_{j=0}^d X_{t-j}\gamma_j} \text{Bias}(\hat{p}_t^\pi) + p_t \left[\frac{\sum_{k=0}^d X_{t-k}\pi_k}{\sum_{j=0}^d X_{t-j}\gamma_j} - 1 \right]
\end{aligned}$$

□

B Alternative data sources

B.1 Retrospective deaths

JHU presented daily deaths in real time, aligned by the date they were reported. In contrast, the National Center for Health Statistics (NCHS) provided weekly totals for deaths aligned by occurrence, and were not

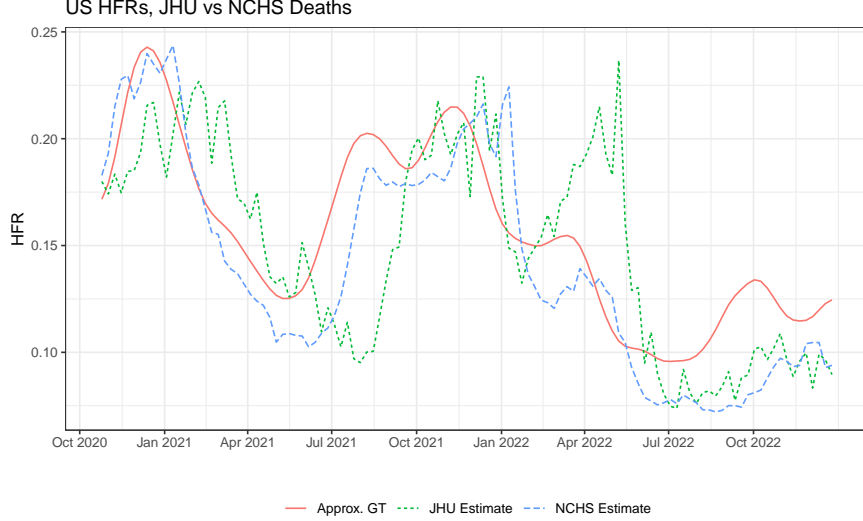


Figure 5: Real-time Lagged Ratios, JHU vs NCHS deaths. Seven-day smoothing with 19- and 11-day lags, respectively.

available in real time. Thus, delay distributions with NCHS deaths have a lighter tail.

Figure 5 shows this minor change has a significant effect on the bias. It compares the real-time lagged ratios (2) with deaths sourced from JHU and NCHS. JHU is much more biased during the variant periods discussed. For example, NCHS only rises from 12% to 14% as Omicron falls, far below JHU’s surge above 25%. As analyzed in Section 2.3, JHU’s heavier-tailed delay distribution inflates the influence of dates with higher HFRs than the present.

B.2 Alternative ground truth

We considered two retrospective approaches to approximate the ground truth national HFRs over time. The first approach took lagged ratios with aggregate deaths from NCHS. NCHS is a better resource than JHU because it uses death counts from the date they actually occurred, not merely reported. In addition, we take a forward-looking ratio, which is retrospective insofar as it uses data after time t to estimate the HFR.

$$\hat{p}_t^{\text{LaggedRetro}} = \frac{Y_{t+L}}{X_t} \quad (8)$$

The second approach computed a single HFR for each major variant, then mixing by the proportions of variants in circulation. Formally, let \hat{p}_j approximate the HFR of variant j ; let v_t^j be its proportion of cases at time t , where $\sum_j v_t^j = 1 \forall t$. The HFR estimate is

$$\hat{p}_t^{\text{Var}} = \sum_j v_t^j \hat{p}_j.$$

Each variant’s HFR \hat{p}_j was defined as the ratio of total NCHS deaths and HHS hospitalizations during the period where it accounted for over 50% of activate cases. The case proportions v_t^j were obtained from [covid.cmu.edu](https://covid.cmu.edu/covidhub/). To ensure estimates were reasonable, we only considered the 4 largest variants: The original strain, Alpha, Delta, and Omicron. Because Omicron began with an enormous surge that quickly subsided, we split it into early and late periods at April 1, 2022, following (Adjai, Stacey and Hong, Kai and Molinari, Noelle-Angelique M and Bull-Ottersson, Lara and Ajani, Umed A and Gundlapalli, Adi V and Harris, Aaron M and Hsu, Joy and Kadri, Sameer S and Starnes, Jon and Yeoman, Kristin and Boehmer, Tegan K, 2022).

Figure 6 displays the three curves approximating the true HFRs. They have nontrivial differences in magnitude, but move more or less in conjunction. To validate our results, we primarily used the rescaled NCHS HFRs as the least problematic of the three. The retrospective NCHS ratios are subject to statistical

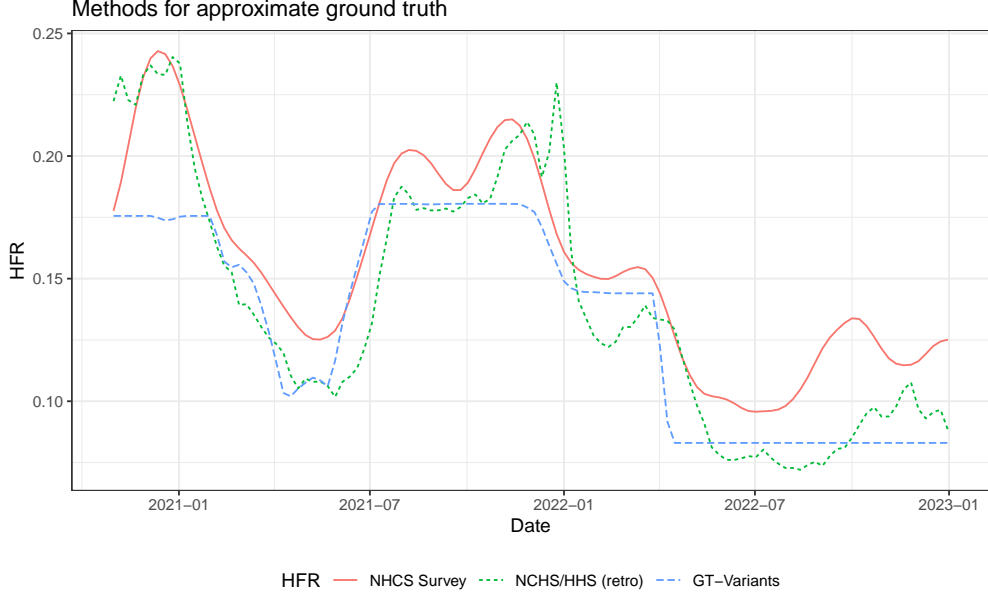


Figure 6: Methods for Retrospective Ground Truth HFRs.

bias, expressed in (7). The variant-based HFRs are flatter, as they do not account for other sources of variability. Therefore, they do not explain for the statistical bias within each variant period, which arises due to changes in the underlying severity rate.

C Robustness checks

C.1 Data Source

The results in Section 3.1 use hospitalization and death counts available in real time. To investigate the sensitivity of our findings, we recomputed the lagged and convolutional ratios, this time using the finalized aggregates. Figure 7a shows the estimates with real-time and finalized counts track very closely to one another. Therefore, the observed bias in 3 cannot be attributed to reporting quirks.

The one period where the curves are significantly different from one another is in March 2022. While the HFRs from finalized counts steadily rise, the real-time estimates sharply fall then immediately bounce back. This sudden drop is due to a brief period in which reported death counts were suddenly too low (Fig. 7b). This is corrected in the finalized counts, hence their smooth HFRs. Removing this artifact further reinforces the bias trends described in Section 2.3.

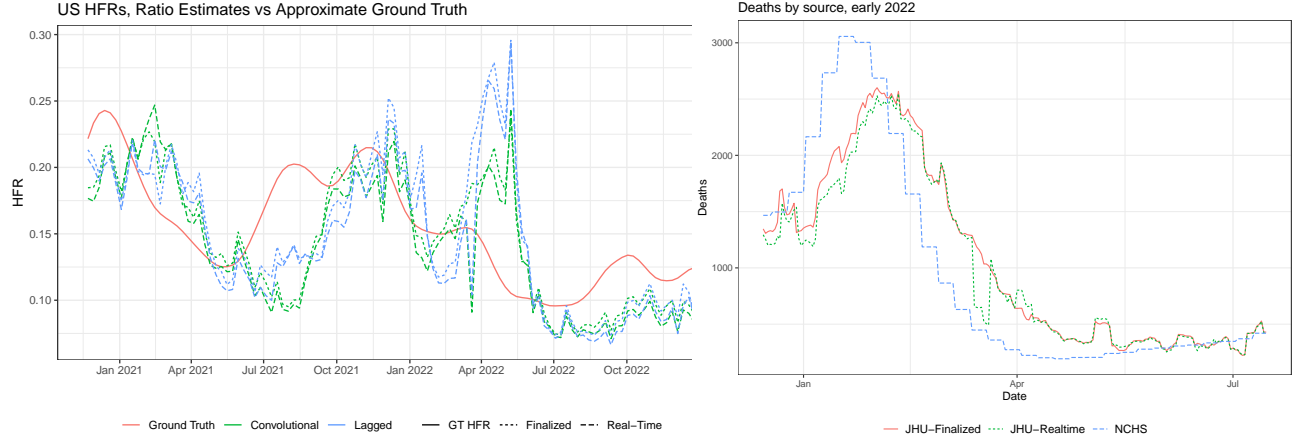
C.2 Hyperparameters

In this section we demonstrate the robustness of our findings against choices of hyperparameters. (All results are with the finalized version of JHU deaths.) First, Figure 8 plots performance over choices of window size parameter. We analyze smoothed versions of the lagged estimator

$$\hat{p}_t^{\ell, W} = \frac{\sum_{s=t-w+1}^t Y_s}{\sum_{s=t-w+1}^t X_{s-\ell}}, \quad (9)$$

as well as the convolutional estimator

$$\hat{p}_t^{\gamma, W} = \frac{\sum_{s=t-w+1}^t Y_s}{\sum_{s=t-w+1}^t \sum_{k=0}^d X_{s-\ell-k} \gamma_k}. \quad (10)$$



(a) HFR estimates with real-time and finalized counts.

(b) Aggregate deaths by data source and report time.

Figure 7: Convolutional ratio estimates are biased regardless of which delay distribution is selected.

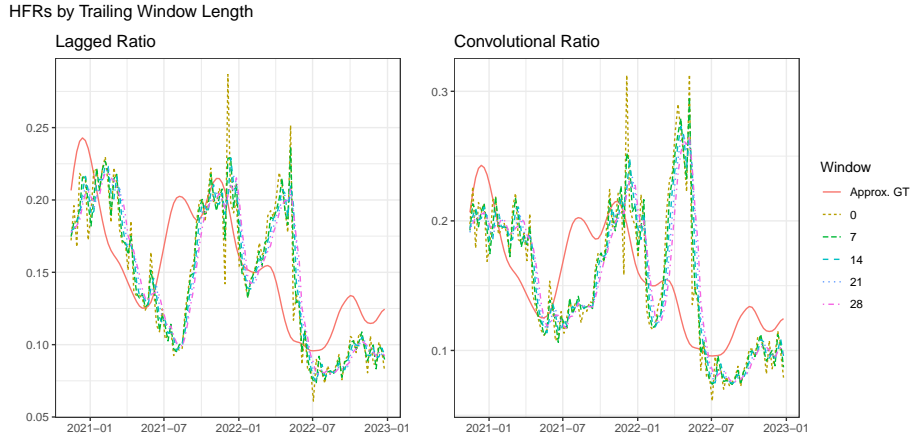


Figure 8: The length of the trailing window bears little impact on the findings.

Results are very similar, indicating the bias does not disappear when smoothing over a longer history.

We next examine the time-to-death hyperparameters: The lag ℓ for the lagged ratio and delay distribution π for the convolutional ratio. Figure 9 displays HFR estimates with lags ranging from 2 to 5 weeks. Unlike the window size, changing this parameter leads to different behavior across lags. Some choices are better than others; a 28-day lag, for example, falls appropriately in winter 2021 and rises less slowly during Delta. However, all are biased to varying degrees, most notably the huge spurious surge in spring 2022.

Figure 10 compares the performance of the convolutional ratio across different choices of delay distribution. We kept the discrete gamma shape for each, but varied the mean and standard deviation. As before, Figure 10a kept the standard deviation to 90% of the mean, per Ward and Johnsen (2021). We also evaluated with a more compact delay distribution in 10b.

All HFR estimates in the figures are significantly biased. Regardless of delay distribution, the ratios are negatively biased during the onset of Delta, and surge after the peak of Omicron. This indicates the bulk of the error is fundamental to the estimator, and cannot be attributed to model misspecification.

Comparing to the approximate ground truth HFRs from NHCS, performance improved slightly with a longer delay distribution than the purported mean of 20 days. Its mean absolute error was 0.031, whereas the delay distribution with mean 28 and standard deviation 25 had a MAE of 0.27. Nevertheless, this difference is relatively small, with the alternative delay distribution still showing similar bias.

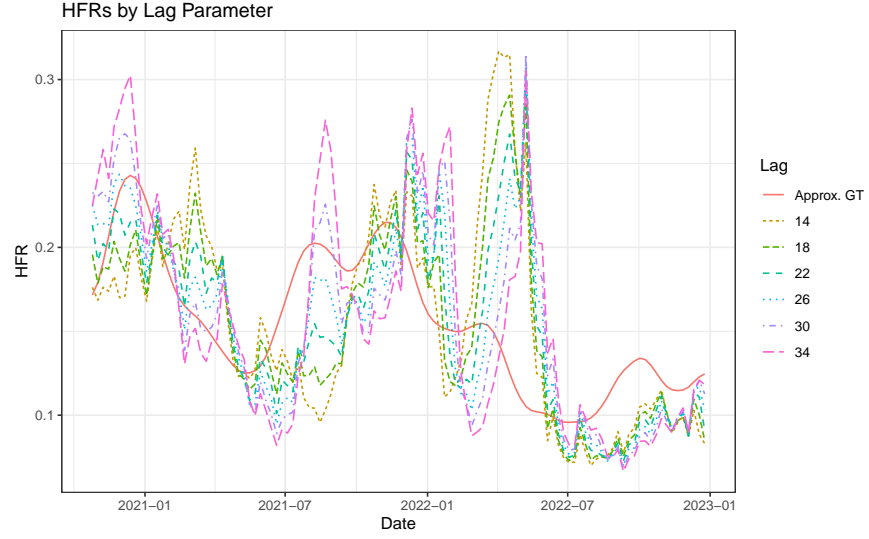
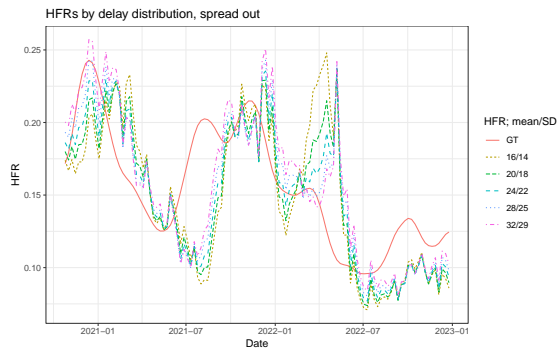
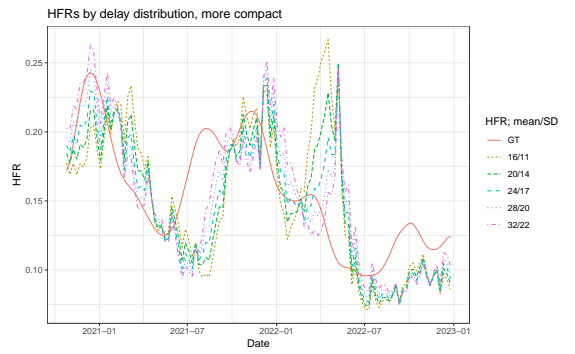


Figure 9: HFRs are biased regardless of what lag parameter is selected.



(a) SD is $0.9 \times \text{mean}$.



(b) SD is $0.7 \times \text{mean}$.

Figure 10: Convolutional ratio estimates are biased regardless of which delay distribution is selected.

C.3 Geography

Next, we repeat our analysis on different geographies, finding similar trends. We repeated our computations on the 6 largest US states with the same lag and delay distribution, with finalized death counts from JHU. Because the NHCS survey was conducted on a subset of hospitals meant to represent the US at large, it may poorly approximate the HFRs for individual states. A better state-level source is the retrospective lagged estimate ((8)) using NCHS deaths. Figure 11 compares this rough ground truth with the real-time estimates. For both NCHS and JHU deaths, we again take the lag that maximizes cross-correlation with hospitalizations; the standard deviation of the delay distribution is 0.9 times the mean.

Several states have similar biases as the US results (Fig. 3a). Ratios in California, Texas, and Florida all are slow to detect the uptick in HFR during Delta; they also spike during Omicron in California, and to a lesser extent Florida. Note these states are the ones with the largest optimal lags, an estimate of the average time to death. As our simulated examples have shown, the shape of the delay distribution is a key factor behind the degree of bias. In contrast, New York, Pennsylvania, and Illinois have mean delays of at most 17. While their HFRs are still biased, they are relatively close to the NCHS curve. This suggests that fatality ratios are generally less trustworthy in states that take longer to report deaths.

D Miscellaneous results

D.1 Alternative lag on simulation

Figure 12 presents the simulation results from Section 3.2 when the lag is the mean of the delay distribution. It clearly indicates the lagged ratio is not markedly better under this lag. This is in step with Figure 9, which demonstrated that its performance on real data are robust to choice of lag.

D.2 Further discussion

In this section, we present examples that further explain the bias. These are more contrived than the ones in Section 2.3, for example using unrealistic delay distributions. Nevertheless, their bias can be simplified to simple analytic formulas, isolating the three contributing factors.

To elucidate the relationship between changing severity rates and the ratio estimators' bias, consider the trivial case where all secondary events occur after exactly ℓ days with no noise. By definition, $\pi_k = \mathbf{1}\{k = \ell\}$, so the convolutional and lagged ratios are both $\hat{p}_t = \frac{X_{t-\ell} p_{t-\ell}}{X_{t-\ell}} = p_{t-\ell}$ presuming both have access to the oracle delay distribution. Figure 13a displays this with the approximate ground truth HFRs from NHCS.

In this case, the bias is the change in the true severity rate $p_{t-\ell} - p_t$. The estimator is unbiased only when the severity rate is stationary. Otherwise, for example, the ratio will be 20% too low if the true severity rate was 20% lower ℓ days ago.

Intuitively, severity rates may be less similar to the present value p_t further back in time. In this simple example, the bias $p_{t-\ell} - p_t$ is generally larger when $\ell = 28$ than $\ell = 14$ (Fig 13a). This expresses the observation that estimates with heavier-tailed delay distributions tend to have more bias.

Section 2.3 claims that changes in primary incidence levels affect the magnitude of bias for the convolutional ratio. Here, we present simple examples that formalize this claim. First assume primary incidence is constant, in which case the convolutional and lagged ratios are equal. The time series factors neatly out of the bias expression Theorem 1:

$$\text{Bias}(\hat{p}_t^\gamma) = \text{Bias}(\hat{p}_t^\ell) = \left(\sum_{k=0}^d \pi_k p_{t-k} \right) - p_t.$$

This is the difference between a weighted average of previous severity rates and the present. Weights for the historical rates are given by the delay distribution, providing further justification for its central role in the bias.

Next, suppose half of the secondary events occur immediately after the primary event ($t = 0$), and the

State-Level True & Estimated HFRs

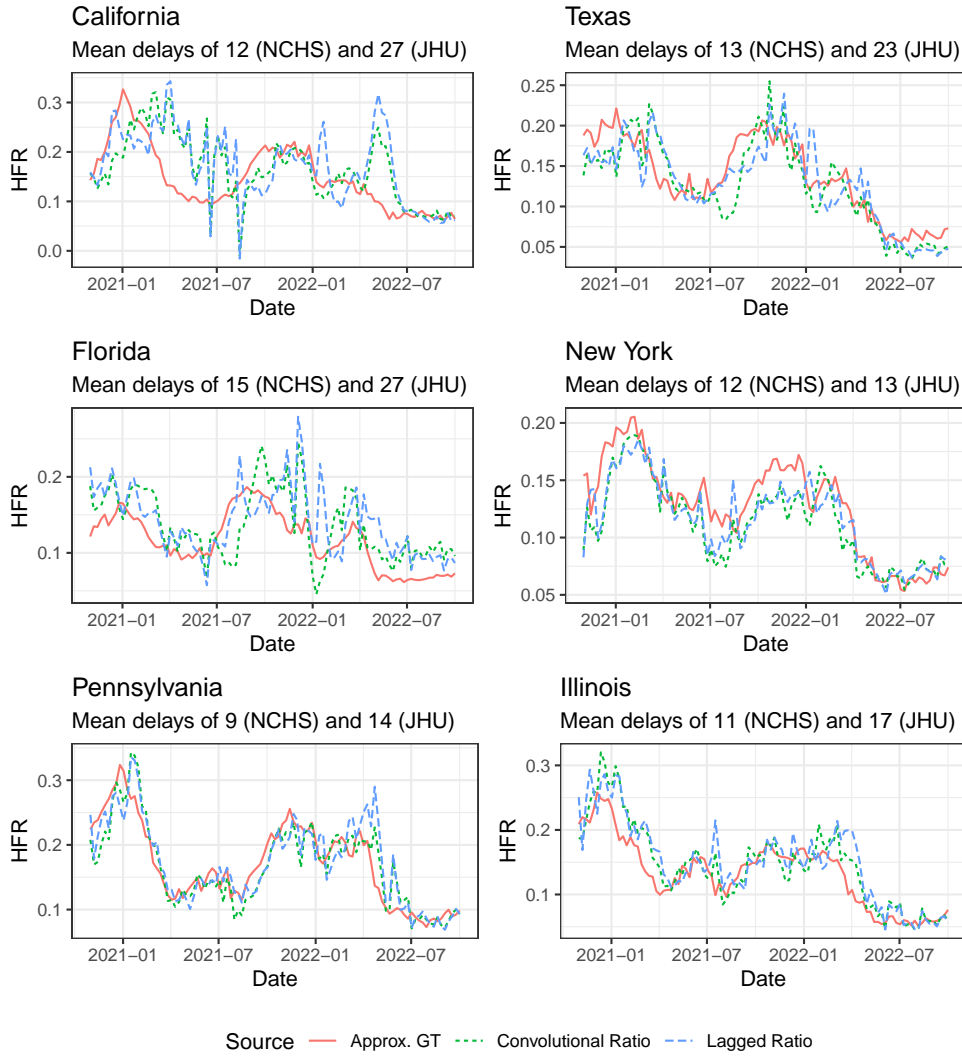


Figure 11: HFRs by individual states. Comparing retrospective estimates with NCHS against real-time estimates with JHU.

HFRs, Simulated Deaths

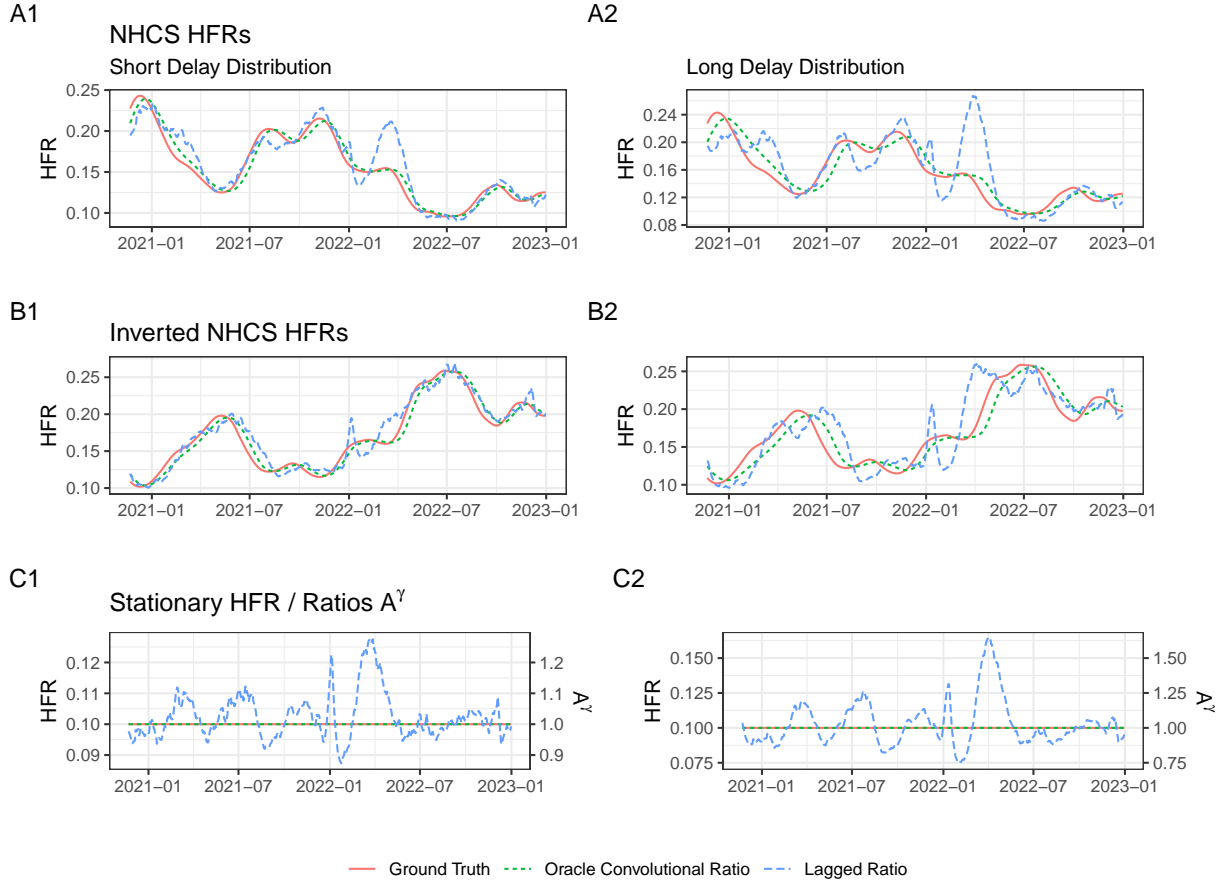
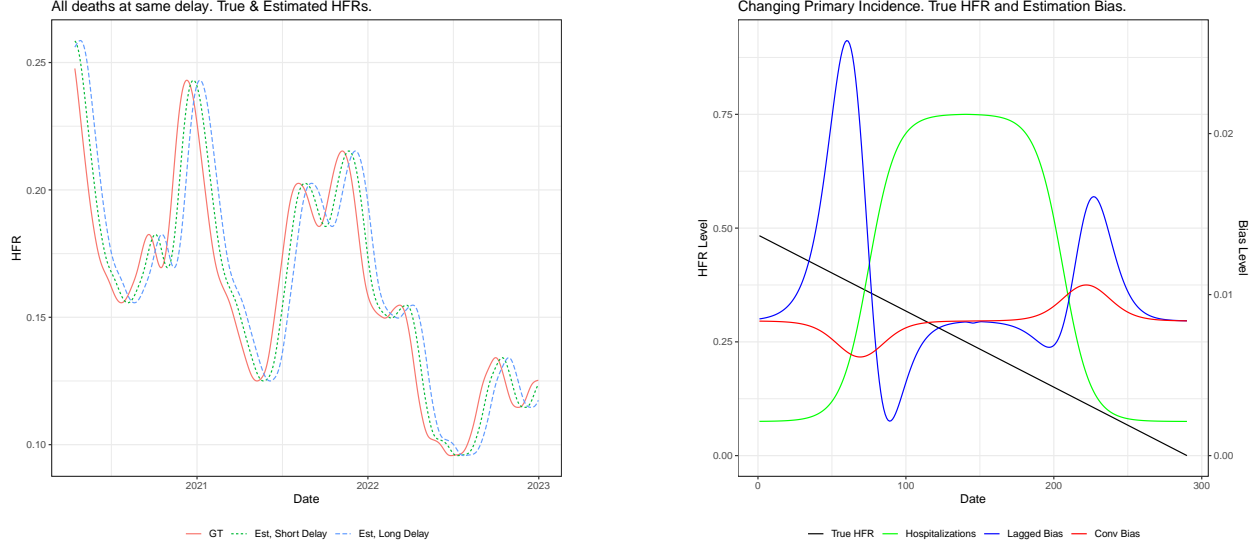


Figure 12: Lag chosen as mean of delay distribution. Same simulated data as Section 3.2.



(a) All deaths after ℓ days. HFR ratios equivalent; plotting delays of $\ell = 14$ and 28 days.

(b) Changing primary incidence. Plotting bias of lagged and convolutional ratios.

Figure 13: Toy examples of biased severity rates.

other half after d days. Further assume $p_{t-d} \neq p_t$, so there is some degree of bias. Then

$$\begin{aligned}
 |\text{Bias}(\hat{p}_t)| &= \frac{\frac{1}{2}|X_t(p_t - p_t) + X_{t-d}(p_{t-d} - p_t)|}{\frac{1}{2}(X_t + X_{t-d})} \\
 &= \frac{X_{t-d}|p_{t-d} - p_t|}{X_{t-d}(1 + \frac{X_t}{X_{t-d}})} = \frac{|p_{t-d} - p_t|}{1 + \frac{X_t}{X_{t-d}}}
 \end{aligned}$$

The absolute bias is monotonically decreasing in $\frac{X_t}{X_{t-d}}$, the proportion change in primary incidence. Rising primary incidence ($\frac{X_t}{X_{t-d}} > 1$) yields less bias, while falling levels yield more.

Figure 13b displays this setting. Hospitalizations are defined as $X = \sigma(s) * 9000 + 1000$, where σ is the sigmoid function and s takes 300 evenly spaced steps from -9 to 7. The true HFRs fall from 0.5 to 0 over the same number of even steps. Indeed, the convolutional ratio's bias dips as hospitalizations rise, and rises as they fall.

When daily hospitalizations approach a constant level, the two estimators become the same ratio, so their biases converge. During periods of change, however, the lagged estimator has different bias. It first moves in the opposite direction as the convolutional bias, before oscillating in the same direction. Moreover, the magnitude of the lagged is far greater.

Since the convolutional ratio uses the true delay distribution, it has oracle bias in Theorem 2. The lagged ratio has lag at the mean of the delay distribution, $\ell = \frac{d}{2}$. Its behavior can be explained by the ratio $A_t^\ell = \frac{X_{t-2\ell} + X_t}{2X_{t-\ell}}$. As hospitalizations begin to steeply rise, $X_{t-2\ell}$ and $X_{t-\ell}$ are similar, but $X_t > X_{t-\ell}$. Therefore $A_t^\ell > 1$, inflating the positive oracle bias term and adding misspecification bias. As hospitalizations level out near the top, $A_t^\ell < 1$, hence the bias falling lower. The opposite pattern occurs as hospitalizations fall.