

WORKING PAPER · NO. 2021-103

Using Household Rosters from Survey Data to Estimate All-cause Mortality during COVID in India

Anup Malani and Sabareesh Ramachandran

DECEMBER 2021

Using household rosters from survey data to estimate all-cause excess death rates during the COVID pandemic in India

Anup Malani & Sabareesh Ramachandran*

December 2021

Abstract

Official statistics on deaths due to COVID undercount deaths due to lack of testing. In developed countries, death registries have been used to measure total excess death during the pandemic. However, very few developing countries have even partial death registries or the capacity to register deaths during a pandemic. In this paper we estimate excess deaths in India using the member roster of a large and representative household panel survey. We estimate roughly 6.3 million excess deaths during the pandemic through August 2021. We cannot demonstrate causality between COVID and deaths but the timing and age structure of deaths is consistent with the COVID pandemic and excess deaths are positively correlated with reported infections. Finally, we find that excess deaths were higher among higher-income persons and were negatively associated with mobility. The methods in this paper can be used in countries with a household panel to measure health-related demographic indicators.

*Malani: University of Chicago Law School and NBER; Ramachandran: UC San Diego. Malani acknowledges funding from the Becker Friedman Institute at the University of Chicago to purchase a subscription to the Consumer Pyramids Household Survey and the support of the Barbara J. and B. Mark Fried Fund at the University of Chicago Law School. We thank Mahesh Vyas, Kaushik Krishnan, Chinmay Tumbe, Shamika Ravi, Rukmini S, Prabhat Jha, Arvind Subramanian, Justin Sandefur, Abhishek Anand, Anmol Somanchi and seminar participants at the CMIE weekly webinar for helpful comments. We thank Bartek Woda and Satej Soman for help with COVID case and deaths data and the anonymous data scientists at www.covid19india.org for scraping data on COVID cases and deaths.

Introduction

COVID is the largest global pandemic since the 1918 flu. According to official reports, over 195 million people have been infected and 4 million have died. Even these remarkable numbers, however, may be an undercount, especially in developing countries. Serological surveys suggest that 20-100 times more people have been infected and have antibodies than have been tested and counted in official reports (e.g., Malani et al., 2020; Mohanan et al., 2021). Likewise, reported deaths may be underreported due to low testing rates, low capacity to register deaths, and possible manipulation of death records (e.g., Rukmini S, 2021).

It is important to accurately measure the number of people infected and dead due to the pandemic in order to estimate the real impact of the pandemic on people's lives. Accurate measurement is also necessary for estimating parameters, such as the infection mortality rate, the impact of non-pharmaceutical interventions, and efficacy of vaccine campaigns, that guide pandemic response.

Total excess deaths during the pandemic is one alternate measure of the mortality risk from the pandemic. In developed countries, data from accurate death registries have helped calculate excess deaths during COVID (e.g. Woolf et al., 2021). However, less than 65% of countries in Asia and 20% countries in Africa have death rates from even partial registries available (United Nations, 2021). Even in countries with a partial registry, the death registration process itself may be affected by the lockdown during the pandemic resulting in error in the excess deaths estimate. Further, governments may have chosen to not report all COVID related deaths either to deflect blame away from them or to prevent people from panicking (e.g., Rukmini S, 2021).

India is one of the hardest hit countries. Officially, India is ranked second globally in number of infections with over 30 million and third in deaths with nearly 500 thousand (Johns Hopkins University & Medicine, Johns Hopkins University & Medicine). These official COVID-death numbers are widely thought to be undercounts (Gamio and Glanz, 2021). However, excess-death estimates based on India's death registries also likely undercount deaths: only 92% of all deaths were recorded by the registries in 2019 (India Ministry of Home Affairs, 2021).

To address this data gap, we use an alternative source of deaths data – household roster of a large, panel data set – to estimate all-cause excess mortality in India during the pandemic. The data set is the Consumer Pyramids Household Survey (CPHS). Its nationally-representative sample includes roughly 174,000 households with roughly 870,000 current members. The survey is conducted on the same households every 4 months, with a representative quarter of the sample

surveyed each month. The survey keeps a roster of all current and past household members and provides reasons for attrition, including death. We count these deaths before COVID to estimate a baseline death rate, and during COVID to calculate excess deaths during the pandemic. An important feature of our data is that it is private and measures death incidentally. This means it is immune to political censorship and is unlikely to have investigator-side bias with respect to death reporting.

In our preferred estimates, the COVID pandemic is associated with 6.3 million excess deaths, over 13 times the number of COVID deaths reported¹. Excess deaths peak in the same months as infections peaked during the two waves that struck India (September 2020 for wave 1 and April-May 2021 for wave 2). Moreover, we find that the second wave experienced significantly greater deaths than the first wave. Although we do not find statistically significant differences in mortality by sex or urban location, we do find that the age-pattern of deaths is COVID-like: deaths rise significantly relative to baseline for those over 60, but decline somewhat for those under 40. We also find a significant correlation between excess deaths in a district and confirmed cases in that district. Incidentally, the excess deaths are higher in families with a higher per-capita income.

Our use of household rosters from a survey to estimate health-related demographic parameters is possible because the data set we employ is representative, large, and repeatedly surveys the same households. However, the use of rosters in this manner has shortcomings that we must address. The main problem is that the survey measures whether a death occurred since the last time the household was surveyed, but does not measure exactly when the death occurred.² We primarily tackle this by restricting the sample to individuals who are observed in consecutive rounds and attribute deaths to the median month between the current and last completed survey. This interprets death rate reported in month t as a moving average of death rates from months $t - 3$ to t . We discuss other methodological issues with using household rosters to track demographics in the appendix.

We compare our estimates with estimates of excess deaths using the Civil Registry System (CRS) data, the official death registries, from 12 states Banaji and Gupta (2021). Our estimates of excess deaths is somewhat larger than the estimate from CRS data. This may be due to the incompleteness of the India's death registries (Deshmukh et al., 2021). We also compare our estimates to those

¹The officially reported number of deaths till 31 August 2021 was 458,470.

²It does not measure why the death occurred either. Therefore, it only allows us to measure excess deaths. In a separate project, we are conducting verbal autopsies on all reported deaths in the survey during 2019-2021 to determine which deaths were plausibly due to COVID.

from the US. Whereas we report a 29% increase in death rates during the pandemic, the US reports a roughly 22% Woolf et al. (2021) increase in excess deaths. Our preferred estimates of deaths are higher than US estimates, though that could be because India has a worse health care system Das et al. (2012).

Our main contribution is to provide novel estimates from India to a growing literature on excess deaths from COVID. Unlike studies from countries that have reliable death registries (Rossen et al., 2020; Woolf et al., 2021; Kontopantelis et al., 2021), it examines a country with unreliable registries. To address the problem of incomplete registries, we employ a large, representative survey (CPHS) that is independent of political influence and allows estimation of heterogeneity in death rates.

Alternative estimates from India employ data on registered deaths but scale them up based on their degree of undercounting (Anand et al., 2021; Deshmukh et al., 2021). However, these estimates are only available for a third of India’s 29 states. Deshmukh et al. (2021) also provides national estimates using other representative surveys. The main advantage of using CPHS over these other surveys is that CPHS has better temporal coverage and tremendous detail on the deceased, providing opportunities to explore whether excess deaths have “COVID-like” features and heterogeneity in death rates. A third approach is to apply estimates of infection fatality rates outside of India to estimates of infection rates in India Anand et al. (2021). The problem with this approach is that India may not have the same infection fatality rates as other countries, just as it does not have the same rates of death from other diseases. Moreover, there are conflicting estimate of seroprevalence in the same place due to antibody decline and many locations lack any seroprevalence estimates. So infection rates have wide confidence bars. One other, contemporaneous paper employs CPHS to estimate excess deaths (Anand et al., 2021). We explore some of the data problems with CPHS a bit more than that paper. To be fair, our estimates are higher than, but not tremendously out of line with available excess-death estimates from other sources.

A second contribution is to show how best to use rosters from household surveys such as CPHS to measure items, like death, migration and marriage, that are implicitly measured by household rosters, in India. To some extent, the problems associated with using CPHS to measure roster-events, especially the timing of these events, are also problem for measuring roster-events in surveys other than CPHS and outside India. Thus, our methods for addressing that may be relevant for counting roster-events from other surveys.

Our analysis has limitations. First, it does not estimate deaths from COVID: it estimates excess deaths during COVID relative to the number of deaths during a control period (e.g., 2019),

before COVID. The deaths could be due to policies cause by COVID, such as lockdowns, or behavior responses to COVID. Second, our estimates depend on our assumptions about pre-pandemic, baseline death rates. Death rates jump from 2018 to 2019, in both CPHS and other data. This is unlikely to be a pre-trend because there is no jump from 2017-2018; indeed, death rates decline from 2015-2018, as is typical over time. Our preferred estimates use 2019 death rates as a benchmark to side step the question of why death rates jumped in 2019. However, if we use 2015-2019 as a baseline, our estimated excess-death rate jumps because 2015-2018 deaths are lower than 2019.

1 Background

According the Global Burden of Disease (GBD) project (Appendix Figure A5) India had a death rate of roughly 0.7% (7 deaths per 1000 persons per year) in 2019, approximately 9.5 million deaths in a population of 1.4 billion (Vos et al., 2020). GBD shows an uptick in the death rate from 2018 to 2019, a pattern also evident—though more pronounced—in our CPHS data.

SARS-CoV-2 hit India in two waves (Figure A6). The first cases were reported on 27 January 2020 (Andrews et al., 2020). The first wave peaked in September 2020, with almost 100,000 confirmed cases and 1,000 deaths daily. The second wave peaked in April 2021, with roughly 400,000 confirmed cases and 4,000 deaths daily (www.covid19india.org, 2021).

India imposed a national lockdown from 24 March to 1 June 2020, well before the wave 1. Google mobility statistics show mobility fell 40% relative to January 2020 levels during that lockdown. After that, lockdowns were local and driven by states. But by the peak of wave 1, mobility had returned to about 15% below January 2020 levels. There were local lockdowns and a reduction of mobility during wave 2, but the decline was not as severe as during wave 1.

Official numbers on cases and deaths should be taken with a grain of salt. Confirmed cases undercount actual infections, at different rates over time. Perhaps 90% of cases were asymptomatic and unlikely to be tested (Waghmare et al., 2021). Per-capita testing rates in India were low relative to developed countries. Testing rates increased dramatically from wave 1 to 2, so the higher case counts may partly be due to testing not cases.

Reported COVID death rates may also be lower than true death rates. First, many deaths in India, especially those outside the hospital setting, are not officially recorded. Second, not all dying individuals are tested for COVID. Further, even individuals dying after a positive COVID test are sometimes recorded as a non-COVID death because they have co-morbidities that could have been

the cause of death. While some such deaths are not causally COVID deaths, some of them may be but are missed.

Because many COVID-attributable deaths are not labeled COVID deaths, researchers have examined all-cause mortality to gauge the impact of the pandemic. Typically the level or projected trend of deaths pre-pandemic is compared to the level of deaths during the pandemic. In India, all-cause mortality is recorded by each state's Civil Registry System (CRS). The main alternative is the Sample Registration System (SRS), which calculates death rates based on a representative, 1% sample of the population.

Each source has its problems. The CRS has three problems. One, not all deaths are registered. For example, in 2017, among big states, the reporting rate was 63.5% in Jammu & Kashmir and 76.4% in Bihar (Rao and Gupta, 2020). Two, while reporting rates are improving over time, this trend complicates estimation. An increase in death rates could be due to better reporting or to an actual increase. Three, CRS reports with delay Ravi (2021). For example, only 14 (of 28) states have CRS data currently available (Rukmini S, 2021).

The SRS also has problems. First, while the CRS is delayed a few months, the SRS is typically delayed 2 years. We may not get SRS estimates of COVID-period deaths until 2023. Second, even the SRS misses about 12% of deaths Gerland (2014). Third, CRS and SRS can diverge. In 2017, the ratio of CRS to SRS deaths range from 38% in Uttar Pradesh to 124% in Tamil Nadu Rao and Gupta (2020). The CRS number can be higher than the SRS number not only because SRS may be an underestimate, but because CRS will report the death of a resident of one state in another state if they went to that other state for medical care and died there Ravi (2021).

Even if one can measure excess deaths, not all are linked to COVID. The death rate due to COVID is typically captured by two epidemiological parameters. The case fatality rate (CFR) is the number of confirmed deaths divided by the number of confirmed COVID cases. This is not a very useful statistic. Both numerator and denominator are undercounted. Moreover, CFR may reflect testing rates and selection into testing as much as harm from the disease. A better alternative is the infection fatality rate (IFR), i.e., COVID deaths divided by all COVID infections (rather than just confirmed infections).

Initial efforts to calculate the IFR used serological studies to estimate the denominator. However, they used official death counts as the numerator Malani et al. (2020); Mohanan et al. (2021); Malani et al. (2021). Because those counts are also underestimates, correcting only the denominator likely led to an underestimate of the IFR. We will not be able to correct that in this paper, as

all-cause excess deaths may include deaths not directly related to COVID. However, it does provide some insight into how off prior IFR estimates might be.

2 Methods

2.1 Data

2.1.1 Consumer Pyramids Household Survey

Our primary data source is the Centre for Monitoring the Indian Economy’s Consumer Pyramids Household Survey (CPHS), a large, representative³, panel survey of Indian households. The sample is based off the 2011 Indian Census and representative at the level of strata defined as homogeneous regions \times urban status and at the national level. A homogeneous region is a cluster of similar districts within a state.⁴ Sample households are visited every 4 months, with each 4-month period called a round. However, a nationally-representative subsample of households are sampled each month. CPHS started in January 2014 and the latest data we could access are from October 2021.⁵

Although the purpose of the CPHS is to measure household economic characteristics, it also maintains a meticulous household roster. The roster records whether there is a death in the household since the last time a household was surveyed, typically 4 months earlier. CPHS also provides data on the demographics and income of each household member, location at the district level, and whether the household resides in a rural area, defined as a village in the 2011 Indian Census. We use these data to explore heterogeneity of death rates.

2.1.2 COVID cases and mobility data

We obtain district-by-day level data on confirmed cases from www.covid19india.org. We obtain estimates of daily infections by scaling up confirmed-case curves with estimates from a serological survey, as we explain in Appendix Section C.

We obtain mobility data from Google’s Community Mobility Reports Google (2021). The units are percent relative to a baseline that is the median value for the corresponding day of week during the 5-week period Jan 3–Feb 6, 2020. Google reports 6 measures of mobility based on location; we

³CPHS has been criticized for not being representative of poor populations (Dreze and Somanchi, 2021). We address this issue in the appendix.

⁴The sampling method is explained in the appendix.

⁵Whereas data through August 2021 are in the People of India file of the CPHS, the September and October 2021 roster is obtained from June 2021 income file, which is released in November 2021.

take an average of the 5 measures other than mobility at home because home mobility rises during the pandemic.

Because our death data are reported monthly, we average daily cases, infections and mobility over each month. We do not have case and mobility data prior to February 2020. Therefore, we assume cases and infections are 0 and that mobility is 100% before that date.

2.2 Data issues with CPHS

Using survey data rather than death registration to measure mortality rates raises a number of data cleaning problems. First, survey response rates fell during the pandemic, particularly during India’s lockdown. Second, there may be selection bias in non-response. Specifically, non-response may be a function of whether a household experienced a death. Third, there appears to be a level jump in the death rate in 2019, prior to the pandemic. Fourth, there CPHS has been criticized for not being representative of poor populations (Dreze and Somanchi, 2021). Finally, the CPHS does not report the precise date death occurred. We address the first 4 concerns in the Appendix and address the timing issue here.

The date on which deaths are “observed” in CPHS is not necessarily the date the death occurred. Households do not report deaths to the CPHS *when* those deaths occur, but some time later when the household is surveyed. Nor does CPHS ask when deaths occurred. So, if a household answered two surveys in a row 4 months apart and reported a death in the second survey, all we know is the death occurred during the intervening 4 months. If the household skipped n surveys between surveys they answered, then a death reported in the last survey occurred sometime in the $4(n + 1)$ months in between responses.

Our preferred solution is to reallocate death to the midpoint between answered surveys⁶. So if there is a gap of k months between surveys the household answered, then a death reported on month t is re-allocated to month $t - (k/2)$. The advantage of this solution is that it is simple. The disadvantage is that it gets the timing of deaths a bit off in the way a moving average would because it smooths out the jump in rates, in part to periods before and in part to periods after the jump.

⁶We explore a second solution in the appendix: estimating the death rate by asking the question, how much would the true death rate have to have changed for the observed death rate to have changed as much as it did since the last month.

2.3 Definition of pandemic period

We define the pandemic period as February 2020 to the present because India’s first confirmed cases are on 27 January 2020 Andrews et al. (2020); the first wave as February - December 2020; and the second wave as January 2021 to the present. The CPHS data are coded at the monthly level, so we cannot more finely define the pandemic or waves. In robustness exercises, we vary the start and stop dates +/- 2 months for our preferred estimation strategy.

2.4 Estimating excess deaths

We estimate excess deaths in two steps. First, we predict monthly death rates \hat{y}_{it} in the absence of the pandemic using data from before the pandemic. Second, we regress the difference between an individual indicator for death and predicted individual death rate, $y_{it} - \hat{y}_{it}$, pandemic or pandemic-wave fixed effects. The coefficients on the second regression give us estimates of excess-death rates for various time periods. We explain both steps below.

Predicted death rates. We consider 4 possible counterfactual death rates: 2 parameterizations of the baseline \times 2 baseline periods.

The baseline can be parameterized as either a level or trend. The level baseline is the mean death rate during baseline: $\hat{y}_{it} = \sum_{s \in B} [(1/N_s) \sum_i y_{is}]$, where t indexes months, B is the baseline period, and N_s is the number of people in the sample in period s . The trend baseline is computed in two steps. First, we estimate the trend during the baseline period with the following regression: $y_{it} = \alpha + \beta(t - t_0) + e_{it}$, where the sample includes $t \in B$. Second, we predict y_{it} for t after the baseline period.

Although there are advantages to using a trend as the baseline, our preferred specification will use levels. The demographic literature typically calculates excess deaths using a baseline trend, for two possible reasons. One is to account for population growth. The literature uses data from death registries, which have information on deaths, but not population. Our data, however, include information on births and deaths. We look at a fixed sample of households and add persons when they are born and lose them when they die.⁷ The other reason is to adjust for changes in death rates. Because death rates change slowly (in the absence of disasters), this argument is only important when making long-term projections. Here we focus on projections for a period less than 2 years so changes in baseline death rates are unlikely to be material. Although we will present results for

⁷In theory there could be addition of household and attrition of households. However, those changes are orthogonal to population growth. Moreover, our results stand even if we hold the sample of households constant.

both baseline parameterizations, our preferred specification will use levels because neither argument for trends is strong in our application and using levels is simpler.

The baseline period can either be 2015-2019 or 2019 alone. The former uses more data, but the latter is more recent. Although we will report results with both periods, our preferred specification is using 2019 alone. First, when using a level baseline, using the latest year’s data mitigates some of the error from not using a trend baseline. Second, and more importantly, there is a jump in death rates from 2017 to 2019 in the CPHS. Other data sources (e.g., the CRS and the Global Burden of Disease) also report a jump in 2019, though the magnitude of the jump is larger in CPHS. The jump is surprising because age-conditional mortality rates typically trend down in the long run and there is no known reason for the jump in 2019 across data sources. Using only 2019 as a baseline side-steps this issue.

Excess-death estimates. We estimate a regression of the form:

$$y_{it} - \hat{y}_{it} = \gamma \mathbb{I}(t \in S) + u_{it} \tag{1}$$

using data from January 2019 onwards. Here y_{it} is an indicator for whether individual i who responded in month $t - 2$ was reported dead in month $t + 2$. Our treatment variable is an indicator for some time period S . Our main specification sets that treatment period to be the whole pandemic period ($S = \text{pandemic}$). Some specifications will replace that with an indicator for each wave or interact these pandemic indicators with age or income indicators. Our estimates of excess deaths come from the coefficients on these time and characteristic indicators. Standard errors are clustered at the village/ward \times month level to account for correlation in reporting of deaths within a locality.

3 Results

Raw data from CPHS aggregated to the national level and without adjusting for the timing-of-death suggests a jump in the death rate during the COVID pandemic. Figure 1A shows the sample size between January 2015 and June 2021 and presents death rate as of the date the deaths are reported (not the date they occurred). There is a drop in response rates during India’s lockdown, which forced to CPHS to sample only about half its households – an issue we address in the Appendix. There is a large rise during Wave 1 and smaller rise in Wave 2. The large rise in Wave 1, however, is partly due to delayed reporting of deaths by households that were omitted from the survey during the lockdown. Actual deaths during wave 1 were likely more spread out.

When we clean the data a bit, we obtain a time series that shows a more moderate increase in death rates during COVID (Table 1B). Specifically, we focus on the sample that includes only households that respond in consecutive rounds, so as to reduce imprecision from reassigning the date of death (i.e., keep the moving average at 4 months rather than 8 or longer). We time shift observed deaths back to month $t - (k/2) = t - 2$ for $k = 2$. Finally, we estimated weighted mean death rates by month, where the weights make responding households nationally representative.⁸ We find that death rates drop in April and May 2020, around the time of the national lockdown, but are at or above the 2019 average in other months of the pandemic. There are three spikes during the pandemic. One spike is June 2020, when lockdown is released, another is September 2020, which wave 1 peaks, and the last is in March - May 2021, when wave 2 peaked. These spikes are significantly greater than adjacent months and all but one 2019 month. The Wave 1 and 2 spikes are greater than even the highest 2019 peak.

3.1 Excess-death estimates

The estimated excess-death rate during the pandemic depends on our baseline specification (Table 1). When we define the baseline as the mean death rate during the period 2015-2019 (column 1), our baseline death rate is 0.787%, not far off from the Global Burden of disease estimate. But excess-death rates during the pandemic are 0.543%, which represents a hard-to-believe increase of over 65%. When we use only 2019 as the baseline (column 2), the baseline death rate rises to 1.038%, which is implausibly high. However, our estimate of the excess-death rate during the pandemic is 0.292%, a relative increase that is somewhat higher but in line with the relative increase in excess deaths from the US (Woolf et al., 2021).⁹ Critically (from a policy perspective), the death rate is significantly greater during the second wave than the first wave (column 5).

When our baseline accounts for trends in the pre-pandemic death rate (columns 3-4, our estimates of the baseline death rate fall to more reasonable numbers. Moreover, estimates of the excess-death rate and number also fall because the baseline death rate, due to the trend, would increase in 2020-21 even without the pandemic. Our preferred estimate remains that in column 2 because we do not have good reason to expect material changes in death rates in such a short period. Moreover, in general death rates fall over time and our baseline trend estimates an increase

⁸This means we use both the weight that makes the sample representative and the non-response factor that makes responding households representative of the sample.

⁹Our estimate of the excess-death rate falls a bit to 0.275% or 0.279% (columns 3 and 4) if we include fixed effects for homogeneous region or district, respectively (Appendix Table *tab:main_estimates*).

in death rates.

Our main specification includes only households that respond in consecutive rounds of the survey (roughly 82% of the sample) to mitigate error from our solution to the timing-of-death problem. If we include in our estimation households that skip rounds, we obtain somewhat higher estimates of excess mortality during COVID (0.323 with responders that skip up to 1 round, 0.369 if up to 2 rounds, Appendix Table A3). This seems inconsistent with the nature of non-response bias we discussed in Section B.2. Recall, however, the sample in which responders had higher death rates was restricted to those who responded in round t and $t+8$, about 64% of the sample. If we include those that did not respond at both those times, the consecutive responders actually have lower death rates.

Our estimate of excess deaths falls if we move forward our estimated start date for the pandemic, possibly because it is counting months before any confirmed cases as pandemic months (Table A4). It rises as we move the start date back 2 months, in part because deaths fall during the lockdown, which occurs in April and May 2020, and pushing back the start date moves the low death rate months out of the pandemic period.

3.2 Heterogeneity of death rates

We first explore heterogeneity in excess deaths along lines that would help gauge whether our estimates are credibly picking up the effect of COVID. We then look at other factors that are policy-relevant.

3.2.1 COVID-related factors

Age. Excess deaths follow a COVID-like pattern with respect to age (Figure 2 below and Table A7 in the Appendix). Excess deaths are significantly positive for higher ages, and insignificant and close to zero for lower ones. This right-skew is somewhat greater in wave 2 than in wave 1.

Gender and location. We do not see COVID-like patterns, however, with respect to gender or location (Table A9). Estimated CFR and IFR is greater among males (Green et al., 2021; Pastor-Barriuso et al., 2020), but we do not find significantly greater excess deaths among males. Likewise, prior studies have shown greater infection rates in urban areas Mohanan et al. (2020); Malani et al. (2021), but we do not find significantly greater excess death in cities during the pandemic.

Infections. Excess deaths are positively correlated with confirmed cases and with infection. Table 2 reports the results of an individual-level regression based on (1), except that we replace

the pandemic indicator with confirmed cases or infections. Cases and infections are reported as monthly averages at the district level. Infections are the same as cases but scaled by seroprevalence estimates. We find that deaths are significantly correlated with cases or infections, even when we add controls for monthly average mobility at the district level. This finding increases the credibility of the claim that excess deaths during the pandemic picks up the effect of COVID. However, one should not interpret the coefficients as case fatality rates (CFR) or infection fatality rates (IFRs) as they may capture COVID deaths that were not included in case or infection counts and include non-COVID deaths.

3.2.2 Policy-relevant factors

Mobility. One of the concerns with using excess mortality to measure the harm from the pandemic is that it picks up both the direct effect of COVID infections as well as indirect effect of association behavioral or policy response. For example, it is possible that the pandemic deterred people from hospitals for fear of getting infected and triggered a lockdown that reduced traffic accidents. We find mixed evidence on the indirect effects of the pandemic.

On the one hand, Figure 1, which shows a sharp drop in death rates in April and May 2020, suggests that India’s national lockdown (March 24 - June 1, 2020) was associated with a sharp reduction in deaths. One should be cautious, however, in interpreting this figure because of the timing of death problem. Because we employ our first solution to this problem, the death rate attributed to, say, April are actually reported in June.

On the other hand, deaths are negatively correlated with Google mobility. Table 2 also reports the results of an individual-level regression based on (1), except that we replace the pandemic indicator with Google’s mobility index. Our estimated coefficient in column 5 is that a 10% reduction in mobility was associated with a 0.5% increase in the death rate.

Income. Because CPHS has information on income, we can also compare excess deaths by income. Serological surveys suggest that, in cities, slums were more affected in wave 1 (Malani et al., 2020). News reports suggest that wave 2 disproportionately affected resident that did not live in slums (Khandekar, 2021). To validate these claims, we add indicators of income terciles and the interaction of those income terciles and pandemic or wave indicators to the regression in (1). This evidence suggests that pandemic had a bigger mortality impact on the top tercile (column 1), but that this imbalance was more pronounced in wave 2 (column 2). The first wave affected the middle and highest tercile more than the lowest. The second wave affected only the top tercile

more than the bottom tercile.

4 Discussion

Our preferred estimates imply that there were 6.3 million excess deaths in India during the pandemic, 1.9 million during wave 1 and 4.3 million during wave II. These estimates only include consecutive responders and are similar to estimates if we include households that skip up to one round of survey.

Our preferred estimate of death is roughly 13x as large as the official number of COVID deaths. This does not prove that COVID caused 13x more deaths, but it does suggest official numbers may be a substantial undercount. It is true that all-cause deaths include non-COVID deaths. But these are excess deaths during the pandemic, so it is likely that COVID directly or indirectly (via policy or behavioral change) is related to these deaths. Of course, we cannot demonstrate causation as we do not have a strictly exogenous introduction of COVID or variation in infections.

We benchmark our findings against estimates from the CRS, the official registry of deaths, in the literature. Our preferred estimate is somewhat higher than the estimates of excess deaths in papers that employ CRS data (Anand et al., 2021; Deshmukh et al., 2021). For example, Banaji and Gupta (2021) extrapolate excess deaths from the 12 of the 14 states for which CRS data are presently available. They estimate excess deaths in all of India to lie between 2.8 and 5.2 million excess deaths from April 2020 - June 2021 depending on how one addresses undercounting by the CRS. Our estimates for this same time period are at 5.2 million – at the high end of CRS estimates.

If our measure of excess deaths is assumed to be due to COVID, that disease easily becomes the leading cause of death in India. Prior to the pandemic, the leading causes of death were non-communicable diseases: cardiovascular disease (2.57 million death annually), chronic respiratory diseases (1.16 million annually) and neoplasms or cancers (0.93 million annually). The leading cause of death from communicable disease was respirator illness and tuberculosis (0.86 million annually).

There are three reasons to believe our estimates are in part picking up the direct effect of COVID. First, excess deaths follow a COVID-like age pattern. Second, excess deaths peak when India’s two waves peak. Third, pandemic period deaths are correlated with the amount of infection in a district.

Our analysis certainly has limitations. First, we do not know the cause of death. Therefore, it is difficult to provide whether official COVID death counts are correct or not. We will attempt to

address this in follow-on work that will conduct verbal autopsies on the deaths in CPHS households since 2019.

Second, CPHS data show a big jump in death rates in 2019. This might cast doubt on the validity of CPHS data or suggest that pre-trend that accounts for the mortality increase we observe. The fact that our estimates of excess pandemic-related deaths are consistent with those from other sources in India, suggests that our use of 2019 as a benchmark is valid for estimating that excess deaths. The fact that there was a change in the age pattern of deaths from 2019 to 2020, but not from 2018 to 2019, suggests that changes in 2020 are not a pre-trend.

Third, we have to impute the timing of death because deaths are sometimes reported months after they occur. We offer a solution to the problem that yields death rates that are at the higher end of CRS estimates after the latter are adjusted for undercounting. Moreover, alternatives such as the CRS have their own problems. CRS and SRS undercounted deaths, requiring estimates or assumptions about undercounting rates. Moreover, CRS is not available for all states and SRS will not be reported for several years.

Finally, one might be concerned about non-random response by households. However, unless one includes households that did not respond for 12 months, our estimates of excess deaths do not change substantially even though our sample comprise 95% of our sample.

References

- Altonji, J. G., T. E. Elder, and C. R. Taber (2005). Selection on observed and unobserved variables: Assessing the effectiveness of catholic schools. *Journal of political economy* 113(1), 151–184.
- Anand, A., J. Sandefur, and A. Subramanian (2021). Three new estimates of india’s all-cause excess mortality during the covid-19 pandemic.
- Andrews, M. A., B. Areekal, K. R. Rajesh, J. Krishnan, R. Suryakala, B. Krishnan, C. P. Muraly, and P. V. Santhosh (2020). First confirmed case of covid-19 infection in india: A case report. *The Indian journal of medical research* 151(5), 490–492.
- Banaji, M. and A. Gupta (2021). Estimates of pandemic excess mortality in india based on civil registration data. *medRxiv*.
- Das, J., A. Holla, V. Das, M. Mohanan, D. Tabak, and B. Chan (2012). In urban and rural india, a standardized patient study showed low levels of provider training and huge quality gaps. *Health affairs* 31(12), 2774–2784.
- Deshmukh, Y., W. Suraweera, C. Tumbe, A. Bhowmick, S. Sharma, P. Novosad, S. H. Fu, L. Newcombe, H. Gelband, P. Brown, and P. Jha (2021). Excess mortality in india from june 2020 to june 2021 during the covid pandemic: death registration, health facility deaths, and survey data. *medRxiv*, 2021.07.20.21260872.
- Dreze, J. and A. Somanchi (2021). View: New barometer of india’s economy fails to reflect deprivations of poor households. *The Economic Times June 21* (June 21).
- Gamio, L. and J. Glanz (2021). Just how big could india’s true covid toll be? *The New York Times May 25, 2021*.
- Gerland, P. (2014). Un population division’s methodology in preparing base population for projections: Case study for india. *Asian Population Studies* 10(3), 274–303.
- Google (2021). Covid-19 community mobility report - india. Report.
- Green, M. S., D. Nitzan, N. Schwartz, Y. Niv, and V. Peer (2021). Sex differences in the case-fatality rates for covid-19—a comparison of the age-related differences and consistency over seven countries. *PLOS ONE* 16(4), e0250523.

- India Ministry of Home Affairs, O. o. t. R. G. (2021). Vital statistics of india based on the civil registration system 2019. Report, Vital Statistics Division Civil Registration System Section.
- Johns Hopkins University & Medicine. COVID-19 Dashboard by the Center for Systems Science and Engineering (CSSE) at Johns Hopkins University (JHU). <https://coronavirus.jhu.edu/map.html>. Accessed: 2021-31-01.
- Khandekar, O. (2021). How india’s slums are faring against the second covid19 wave. *Livemint April 29*(April 29).
- Kontopantelis, E., M. A. Mamas, J. Deanfield, M. Asaria, and T. Doran (2021). Excess mortality in england and wales during the first wave of the covid-19 pandemic. *Journal of Epidemiology and Community Health 75*(3), 213.
- Malani, A., S. Ramachandran, V. Tandel, R. Parasa, S. Sudharshini, V. Prakash, Y. Yoganathan, S. Raju, and T. Selvavinayagam (2021). Seroprevalence in tamil nadu in october-november 2020. *MedRxiv*.
- Malani, A., D. Shah, G. Kang, G. N. Lobo, J. Shastri, M. Mohanan, R. Jain, S. Agrawal, S. Juneja, S. Imad, and U. Kolthur-Seetharam (2020). Seroprevalence of sars-cov-2 in slums versus non-slums in mumbai, india. *The Lancet Global Health*.
- Mohanan, M., A. Malani, K. Krishnan, and A. Acharya (2020). Prevalence of covid-19 in rural versus urban areas in a low-income country: Findings from a state-wide study in karnataka, india. *medRxiv*, 2020.11.02.20224782.
- Mohanan, M., A. Malani, K. Krishnan, and A. Acharya (2021). Prevalence of sars-cov-2 in karnataka, india. *JAMA 325*(10), 1001–1003.
- Pastor-Barriuso, R., B. Pérez-Gómez, M. A. Hernán, M. Pérez-Olmeda, R. Yotti, J. Oteo-Iglesias, J. L. Sanmartín, I. León-Gómez, A. Fernández-García, P. Fernández-Navarro, I. Cruz, M. Martín, C. Delgado-Sanz, N. Fernández de Larrea, J. León Paniagua, J. F. Muñoz-Montalvo, F. Blanco, A. Larrauri, and M. Pollán (2020). Infection fatality risk for sars-cov-2 in community dwelling population of spain: nationwide seroepidemiological study. *BMJ 371*, m4509.
- Rao, C. and M. Gupta (2020). The civil registration system is a potentially viable data source for reliable subnational mortality measurement in india. *BMJ global health 5*(8), e002586.

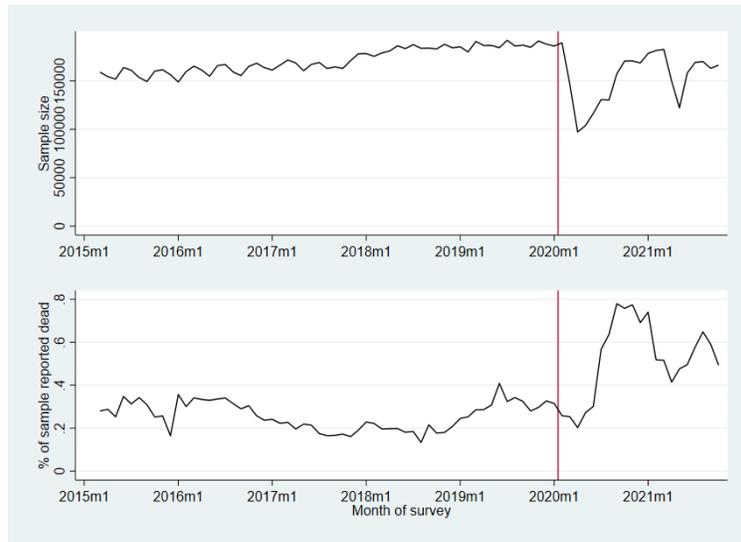
- Ravi, S. (2021). Counting deaths in india is difficult. *Hindustan Times July 14, 2021* (July 14, 2021).
- Rossen, L. M., A. M. Branum, F. B. Ahmad, P. Sutton, and R. N. Anderson (2020). Excess deaths associated with covid-19, by age and race and ethnicity - united states, january 26-october 3, 2020. *MMWR. Morbidity and mortality weekly report 69*(42), 1522–1527.
- Rukmini S (2021, July 6, 2021). Gauging pandemic mortality with civil registration data. *The Hindu*.
- Sharma, H. (2021, July 21, 2021). Two-thirds of indians have covid antibodies, 40 crore still at risk: Icmr. *Indian Express July 21, 2021* (July 21, 2021).
- United Nations, S. D. (2021). Coverage of birth and death registration.
- Vos, T., S. S. Lim, C. Abbafati, K. M. Abbas, M. Abbasi, M. Abbasifard, M. Abbasi-Kangevari, H. Abbastabar, F. Abd-Allah, A. Abdelalim, M. Abdollahi, I. Abdollahpour, H. Abolhassani, V. Aboyans, E. M. Abrams, L. G. Abreu, M. R. M. Abrigo, L. J. Abu-Raddad, A. I. Abushouk, A. Acebedo, I. N. Ackerman, M. Adabi, A. A. Adamu, O. M. Adebayo, V. Adekanmbi, J. D. Adelson, O. O. Adetokunboh, D. Adham, M. Afshari, A. Afshin, E. E. Agardh, G. Agarwal, K. M. Agesa, M. Aghaali, S. M. K. Aghamir, A. Agrawal, T. Ahmad, A. Ahmadi, M. Ahmadi, H. Ahmadi, E. Ahmadpour, T. Y. Akalu, R. O. Akinyemi, T. Akinyemiju, B. Akombi, Z. Al-Aly, K. Alam, N. Alam, S. Alam, T. Alam, T. M. Alanzi, S. B. Albertson, J. E. Alcala-Rabanal, N. M. Alema, M. Ali, S. Ali, G. Alicandro, M. Alijanzadeh, C. Alinia, V. Alipour, S. M. Aljunid, F. Alla, P. Allebeck, A. Almasi-Hashiani, J. Alonso, R. M. Al-Raddadi, K. A. Altirkawi, N. Alvis-Guzman, N. J. Alvis-Zakzuk, S. Amini, M. Amini-Rarani, A. Aminorroaya, F. Amiri, A. M. L. Amit, D. A. Amugsi, G. G. H. Amul, D. Anderlini, C. L. Andrei, T. Andrei, M. Anjomshoa, F. Ansari, I. Ansari, A. Ansari-Moghaddam, C. A. T. Antonio, C. M. Antony, E. Antriyandarti, D. Anvari, R. Anwer, J. Arabloo, M. Arab-Zozani, A. Y. Aravkin, F. Ariani, J. Ärnlöv, K. K. Aryal, A. Arzani, M. Asadi-Aliabadi, A. A. Asadi-Pooya, B. Asghari, C. Ashbaugh, D. D. Atnaflu, et al. (2020). Global burden of 369 diseases and injuries in 204 countries and territories, 1990x2013;2019: a systematic analysis for the global burden of disease study 2019. *The Lancet 396*(10258), 1204–1222.
- Vyas, M. (2021). View: There are practical limitations in cmie’s cphs sampling, but no bias. *The Economic Times June 23* (June 23).

- Waghmare, R., R. Gajbhiye, N. N. Mahajan, D. Modi, S. Mukherjee, and S. D. Mahale (2021). Universal screening identifies asymptomatic carriers of sars-cov-2 among pregnant women in india. *European Journal of Obstetrics, Gynecology, and Reproductive Biology* 256, 503–505.
- Woolf, S. H., D. A. Chapman, R. T. Sabo, and E. B. Zimmerman (2021). Excess deaths from covid-19 and other causes in the us, march 1, 2020, to january 2, 2021. *JAMA* 325(17), 1786–1789.
- www.covid19india.org (2021). Coronavirus outbreak in india. Report.

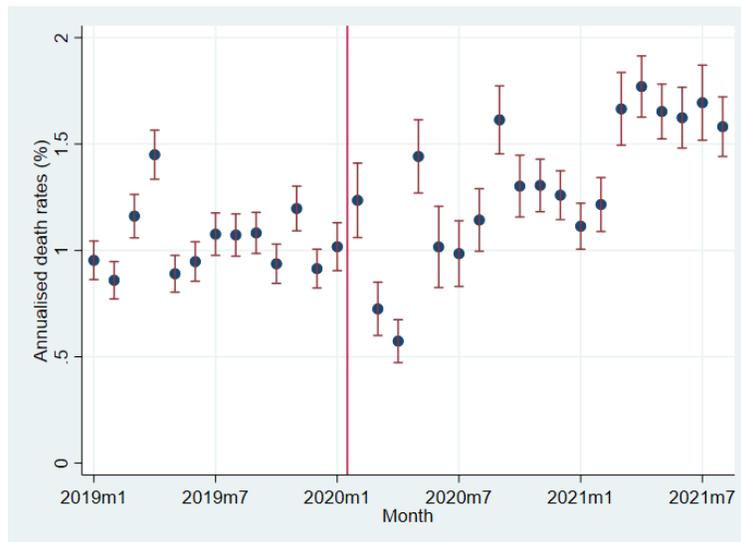
Figures

Figure 1: Sample size, deaths, and death rates.

Panel A: Sample size and households reporting a death from January 2015 - June 2021.

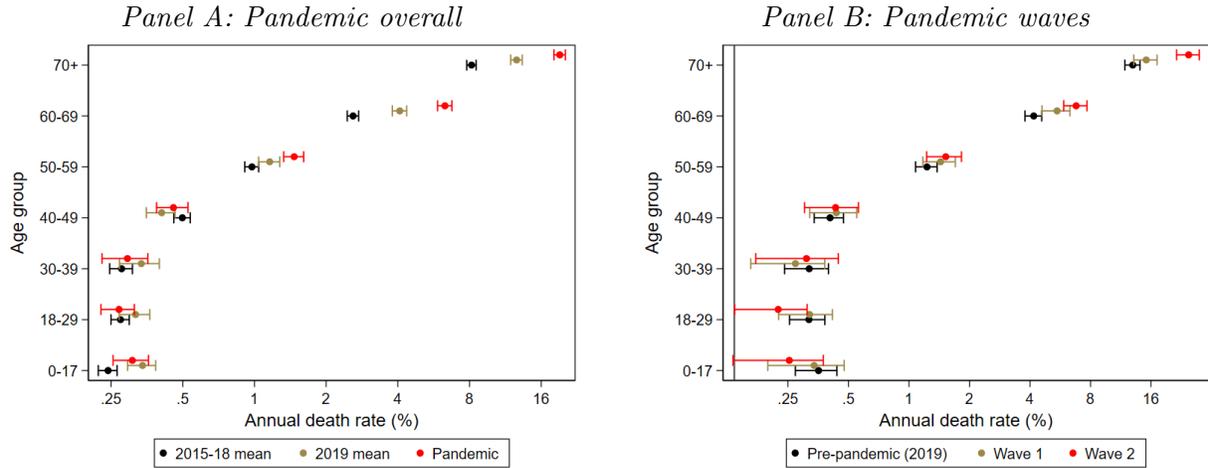


Panel B: Time series of death rates from month fixed effects, January 2019 - July 2021.



Notes: Panel A: The sample includes all responding households regardless of how frequently they respond to a survey. Deaths reported in month t are *not* allocated to a prior month. The y-axis in the lower graph is the proportion of individuals responding in that month who are reported to be dead. The data are not weighted to be representative. Panel B: The sample includes households that respond in consecutive rounds. Observations on individuals are weighted to be nationally representative even with non-response. Each point is the weighted mean death rate in a month and each whisker is the 95% confidence interval on that mean. Both panels: The red line demarcates the start of the pandemic in February 2020.

Figure 2: Age pattern of death rates during the pandemic



Notes. Estimates are from a regression of an indicator for individual deaths on time period indicators: $y_{it} = \sum_s \beta_s \mathbb{I}(t \in s) + w_{it}$. Here y_{it} is an indicator for whether individual i who responded in month $t - 2$ was reported dead in month $t + 2$ and s are indicators for different periods. We only include individuals who reported in the prior survey attempt and the current one. We run separate regressions for respondents in each age group listed in the y-axis. Coefficients on period indicators estimate death rates during those periods for the relevant sample population. In Panel A, we use data from 2015 onwards report coefficients from an indicator for 2019 and for the pandemic period. In Panel B, we use data from 2019 onwards and include indicators for 2019 and the two waves.

Tables

Table 1: COVID death rate prior to the pandemic and the excess-death rate during the pandemic

	(1)	(2)	(3)	(4)
Baseline specification				
Period	2015-19	2019	2015-19	2019
Trend	No	No	Yes	Yes
Deaths during baseline				
Rate (annualised %)	0.787	1.038	0.990	0.872
Excess deaths during pandemic				
Rate (annualised %)	0.543	0.292	0.394	0.244
	(0.0432)	(0.0504)	(0.0417)	(0.0417)
Number (millions)	11.71	6.290	8.501	5.252
	(0.931)	(1.086)	(0.899)	(0.899)
Excess death numbers by waves				
Wave 1 (millions)	5.051	1.914	3.363	1.425
	(0.632)	(0.697)	(0.626)	(0.626)
Wave 2 (millions)	6.617	4.335	5.101	3.789
	(0.618)	(0.670)	(0.603)	(0.603)

Notes. The sample includes households that respond in consecutive rounds. Observations on individuals are weighted to be nationally representative even with non-response. The regression model is (1). in the first two columns we assume a constant counterfactual death rate and in the last two columns we assume the counterfactual death rate has a linear trend. Excess deaths are calculated for the 19 month period from Feb 2020 to Aug 2021.

Table 2: Correlation of death rates with cases, infections, and mobility

	Annualised death rates (%)				
	(1)	(2)	(3)	(4)	(5)
Infections (sero scaled)	0.0195*** (0.00332)	0.0183*** (0.00325)			
Cases			0.00000131*** (0.000000334)	0.00000110** (0.000000337)	
Mobility		-0.000484** (0.000175)		-0.000445** (0.000160)	-0.000749*** (0.000171)
2019 mean	1.120*** (0.0279)	1.082*** (0.0279)	1.151*** (0.0275)	1.116*** (0.0273)	1.126*** (0.0269)
N	2792964	2753588	2783348	2783348	2783348

Notes. Estimates in columns 1 and 2 are from a regression of death rates against sero scaled infections from an SIR model. Estimates in columns 3 and 4 are from a regression of death rates against officially reported COVID cases. We control for mobility in columns 2 and 4. Estimates in column 5 are from a regression of death rates against google mobility. Standard errors clustered at the village/ward \times month level are reported in parentheses. $p < 0.05/0.01/0.001$.

Table 3: Annualised death rates by income groups

	Annualised death rates(%)					
	(1)	(2)	(3)	(4)	(5)	(6)
Pandemic	0.168*		0.311*		0.130	
	(0.0685)		(0.127)		(0.0792)	
Pandemic × 2nd tercile	0.0991		-0.0412		0.134	
	(0.0831)		(0.118)		(0.105)	
Pandemic × 3rd tercile	0.357***		0.240		0.360*	
	(0.101)		(0.132)		(0.140)	
Wave 1		-0.0247		0.0416		-0.0411
		(0.0770)		(0.132)		(0.0900)
Wave 1 × 2nd tercile		0.239*		-0.00842		0.338**
		(0.103)		(0.146)		(0.131)
Wave 1 × 3rd tercile		0.367**		0.154		0.569**
		(0.113)		(0.137)		(0.183)
Wave 2		0.432***		0.644***		0.372***
		(0.0951)		(0.193)		(0.108)
Wave 2 × 2nd tercile		-0.0954		-0.0796		-0.154
		(0.108)		(0.159)		(0.136)
Wave 2 × 3rd tercile		0.330*		0.350		0.0661
		(0.156)		(0.199)		(0.175)
2019 mean	1.160***	1.160***	1.170***	1.170***	1.158***	1.158***
	(0.0466)	(0.0466)	(0.0770)	(0.0770)	(0.0544)	(0.0544)
2nd tercile	-0.172**	-0.172**	-0.150*	-0.150*	-0.184**	-0.184**
	(0.0545)	(0.0545)	(0.0700)	(0.0700)	(0.0677)	(0.0677)
3rd tercile	-0.232***	-0.232***	-0.232***	-0.232***	-0.242**	-0.242**
	(0.0556)	(0.0556)	(0.0688)	(0.0688)	(0.0774)	(0.0774)
Sample	All	All	Urban	Urban	Rural	Rural
N	2969850	2969850	1960145	1960145	1009705	1009705

Notes. Estimates are from a regression model based on 1, with the addition of a income tercile indicator and income tercile indicator interacted with the pandemic or wave indicator. For each individual we calculate the income per capita in 2018. We compute the household's income percentile in their homogeneous region and region type (urban/rural). Households between 33 and 67 percentile are in income tercile 2 and households between 67 and 100 percentile are in income tercile 3. Columns 1 and 2 includes all data, columns 3 and 4 include only urban regions and columns 5 and 6 include only rural regions. Sample includes only consecutive observations and is weighted to be nationally representative. Standard errors clustered at the village/ward × month level are reported in parentheses. $p < 0.05/0.01/0.001$.

Appendix

A Sampling method in CPHS

The country is divided into 99 of these regions. Rural areas are defined as 2011 Census villages. Urban areas are towns and cities. In the rural areas, villages are randomly selected. Within each village 16 households are selected by randomly picking a cluster of homes and then conducting systematic sampling. In the urban areas, towns are divided into substrata based on population. Within each size substrata, towns are randomly selected. Within towns, census enumeration blocks are randomly selected. In each CEB 16 households are selected.

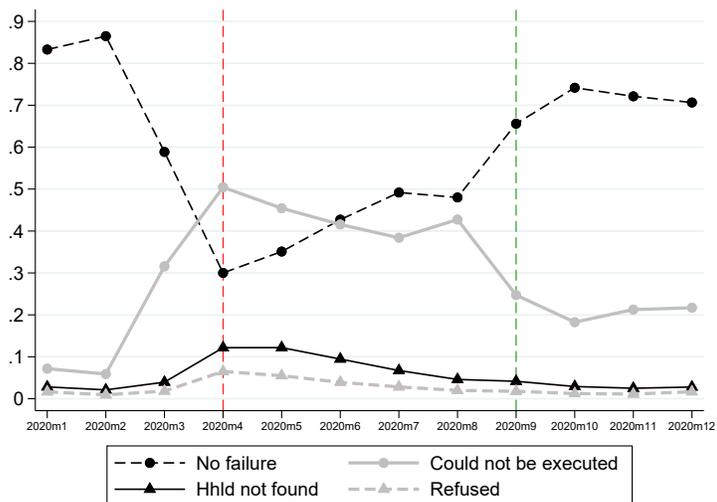
B Additional issues with CPHS mortality data

Recall that using survey data rather than death registration to measure mortality rates raises a number of data cleaning problems. First, survey response rates fell during the pandemic, particularly during India’s lockdown. Second, there may be selection bias in non-response. Specifically, non-response may be a function of whether a household experienced a death. Third, there appears to be a level jump in the death rate in 2019, prior to the pandemic. Fourth, there CPHS has been criticized for not being representative of poor populations (Dreze and Somanchi, 2021). We address these here.

B.1 Low response rate during lockdown

The CPHS experienced a sharp decline in response rates during the lockdown in India. CPHS is ordinarily an in-person survey and the typical per-round, household response rate (responding/sample households) prior to the pandemic was roughly 85%. However, when India’s central government declared a lockdown on March 24, 2020, in person surveys had to cease. CPHS made two changes: they switched to a phone survey and surveyors’ managers, rather than surveyors, conducted the survey to keep up the quality of surveys. Because there are fewer managers than surveyors, CPHS decided only to call a quasi-random, representative subsample of households. The asked managers to pick household phone numbers with only information on strata (defined above) of households and required that the ratio of urban-to-rural households in each homogeneous remain the same as intended pre-pandemic. As a result, response rates fell. Figure A1 shows that the fraction of

Figure A1: CPHS non-execution and non-response rates during 2020.



Notes: Red line indicates first month of phone surveys. Green line indicates month that in-person surveys resumed. The sample includes all households. “Could not be executed” (non-execution) includes both CMIE’s decision not to contact a household and its inability to speak to a household member because, e.g., no one answered the door (“door-lock”). “Household not found” means CMIE attempted to contact the household but surveyors were unable to locate the household.

households that were not contacted rose to roughly 50% of the full sample from April-August 2020 and responding households constituted just 30% of the full sample at the height of the lockdown in April-May 2020, implying a response rate of roughly 60% in April 2020. When CPHS finished its second round in August 2020, it returned to in-person surveys. However, the response rate only rose to 75%. There was also a drop in response rates during wave 2, during which local lockdowns forced local use of telephonic surveys as before.

Low response rates are themselves not an issue because we only look at the proportion of *responding* people who are dead. Also, even if a household does not respond in one round, they may respond in subsequent rounds. In a robustness check we do not restrict to households that respond in consecutive rounds. All reported deaths are accounted for in this estimation.

B.2 Non-random response during COVID

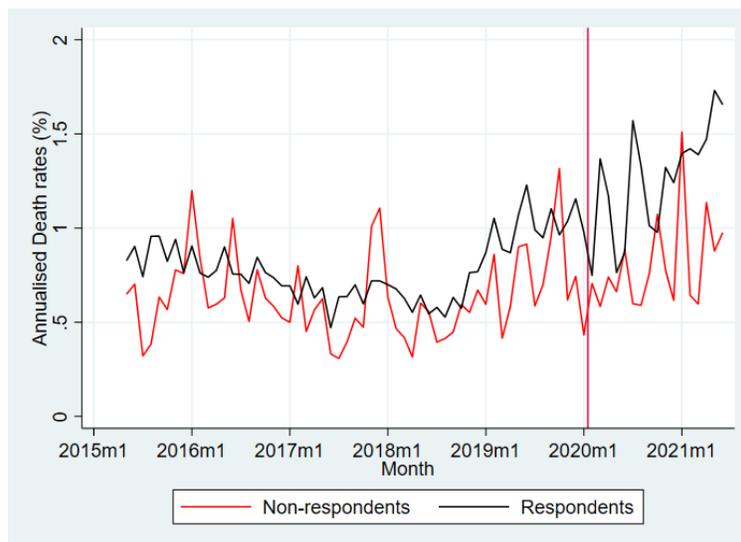
It is possible that households that respond to the CPHS are not representative of the CPHS sample or that the nature of non-random response changed during COVID. The former affects our estimates of baseline death rates unless CPHS’s weights make non-responding households representative even

Table A1: Fit (R^2) from LASSO-generated prediction model for CPHS non-execution and/or non-response in 2019 and 2020 using observables (other than death).

	March-June		Sept-Dec	
	2019	2020	2019	2020
Non-execution & non-response	.004	.008	.007	.007
Non-execution	.005	.009	.010	.008
Non-execution or response	.003	.008	.003	.003

Notes. The sample includes all households in the time period indicated in the column header. The table reports the fit (R^2) from a regression of an indicator for the action in the row label on (a) strata fixed effects (homogeneous region x community type) and (b) covariates selected by LASSO from the set of all variables (excluding death) on the household and its members for the same time period. Observations are at the household level. No weights are included.

Figure A2: Annualized death rates for households that responded in round $t + 1$ and those that did not.



Notes: The sample includes households that responded in round t and $t + 2$. Some of these households also responded in $t + 1$ (the $t + 1$ “respondent” group) and some did not (the “non-respondent” group). This figure plots the annualized death rate in the respondent and non-respondent groups by month of the $t + 1$ survey.

with non-random response.¹⁰ Even if this the former is not true, the latter affects our estimate of the excess-death rate during COVID.

We have mixed evidence about the representativeness of the responding sample. The optimistic view comes from a simple exercise in the vein of Altonji et al. (2005). CPHS has hundreds

¹⁰The former may not affect our estimate of excess deaths during COVID if non-random response is such that the trend of deaths in responding households is similar to the trend in non-responding households, though this is an optimistic assumption.

of variables on each household, including current and lagged responses to income, time use and consumption questions. We estimated a regression of survey response on covariates (other than death) selected via LASSO prior in 2019 and then in 2020. Our estimated R^2 was < 0.01 (Table A1). Of course it could be that survey response is a function of unobservables even conditioned on observables. But given how many observables we have, this seems unlikely. If we make the assumption in, e.g., Altonji et al. (2005) that unobservables have the same explanatory power as observables, then our low bound on the the R^2 from observables implied low R^2 for unobservables.

On the other hand, we do have some evidence that the fact of death affects response rates. This comes from the following exercise. First, we took the set of households that responded in round t and $t + 2$, about 64% of the sample. (Since rounds are 4 months long, this means 8 months apart.) Some of these households also responded in $t + 1$ (the $t + 1$ “respondent” group) and some did not (the “non-respondent” group). (Respondents are 83% of the CPHS subsample that responds at t and $t + 8$.) Second, we compare the number of deaths that occurred between t and $t + 2$ in the respondent group and the non-respondent group. Recall that, even if a household does not respond at $t + 1$, they eventually report their deaths in $t + 2$, so we see all deaths in this period for both groups. Figure A2 plots the annualized death rate in the respondent and non-respondent groups by month of the $t + 1$ survey. Respondent households have slightly more deaths prior to COVID, though in some months non-respondent death rates rise above response ones. But in 2020, the gap widens and the respondent groups death rates almost always appear to be above non-respondent group rates.

Table A2 provides estimates of the sort of bias one would get if one focused only on households responding in consecutive rounds as opposed to on households that at least responded in round t and $t + 2$. The latter are about 64% of the sample. A regression of death rates on a pandemic indicator, a respondent household indicator and the interaction of the two indicators reveals that respondent death rates are 0.276 percentage points higher per annum before the pandemic, and rise 0.182% percentage points further above non-respondents during the pandemic, with each difference being statistically significant (Table A2). Our finding that, of households that respond at t and $t + 8$, households with a death are more likely to respond to a survey does not imply that is true for the full sample. Indeed, as we explain in the next paragraph, the bias is different for the remainder of CPHS that does not respond at t and $t + 8$. So the important take-away is that there could be non-random survey response, not that we can sign it.

There is a solution to non-random response, but it has separate problems. Non-random response

Table A2: Death rate by consecutive survey response status

	Annualised death rates(%)
	(1)
Pandemic	0.0735 (0.0611)
Respondent	0.281*** (0.0517)
Pandemic × Respondent	0.221** (0.0721)
2019 Non-respondent mean	0.732*** (0.0460)
<i>N</i>	2598894

Notes. This table reports the results from regressing an indicator for whether an individual died against an indicator for the duration of the pandemic, an indicator for the response status and the interaction of these two. The sample includes individuals who are observed in month t and month $t+8$, about 64% of the entire sample. The dependent variable for an individual in month t is whether their death status was reported by month $t+8$. Respondent is an indicator for whether the individual also responded in month $t+4$, whether or not they were alive that month. Observations on individuals are weighted to be nationally representative even with non-response. Standard errors clustered at the village/ward × month level are reported in parentheses. */**/** indicates $p < 0.05/0.01/0.001$.

Table A3: excess-death rates for different samples of responders

	Annualised death rates(%)			
	(1)	(2)	(3)	(4)
Pandemic	0.308*** (0.0313)	0.323*** (0.0293)	0.369*** (0.0292)	0.412*** (0.0293)
2019 mean	1.038*** (0.0189)	1.075*** (0.0185)	1.083*** (0.0185)	1.097*** (0.0187)
Sample composition	Responses 4 mo apart	Responses 4/8 mo apart	Responses 4/8/12 mo apart	All
% of full sample	81 %	95 %	98 %	100 %
<i>N</i>	3062171	3560964	3692724	3760155

This table reports the results of the main regression where we regress deaths against a pandemic indicator. All deaths are assigned to the month in the middle of the period when the individual was last surveyed and when they were reported to be dead. The pandemic indicator is set to 1 iff the middle month is after Jan 2020. In the first column we only include household responses in those months for which they responded again in the next round. In the second and third columns we include response if the household responded in at least one of the next two and three rounds respectively. In the fourth column, we include the entire sample. Standard errors clustered at the village/ward × month level are reported in parentheses. */**/** indicates $p < 0.05/0.01/0.001$

can be addressed by using both respondent and non-respondent households in estimate excess deaths. After all, non-respondents at $t + 1$ typically respond at $t + 2$ and this fills in holes in the death data. However, there is loss of precision about the timing of deaths when $t + 1$ non-respondents are included in the analysis: their deaths have to be allocated over 8 months rather than just 4 months. Table A3 presents estimates of baseline death rate in 2019 and excess deaths during COVID as we vary the sample to include more or less non-consecutive respondents. The first column is our main sample of consecutive responders. The second and third allow into the sample households that skip at most 1 and at most 2 rounds of the survey, respectively. The last column includes all households. Adding non-responders increases our estimates of baseline mortality and excess deaths during COVID. The estimates rise rather than fall because none of the columns in Table A3 have the same sample as that in Table A2.

Because we believe that the timing problem is greater than the non-random response problem, we highlight results using consecutive observations in the main text, and report estimates using even non-consecutive observations here in the Appendix.

B.3 Rise in deaths in 2019

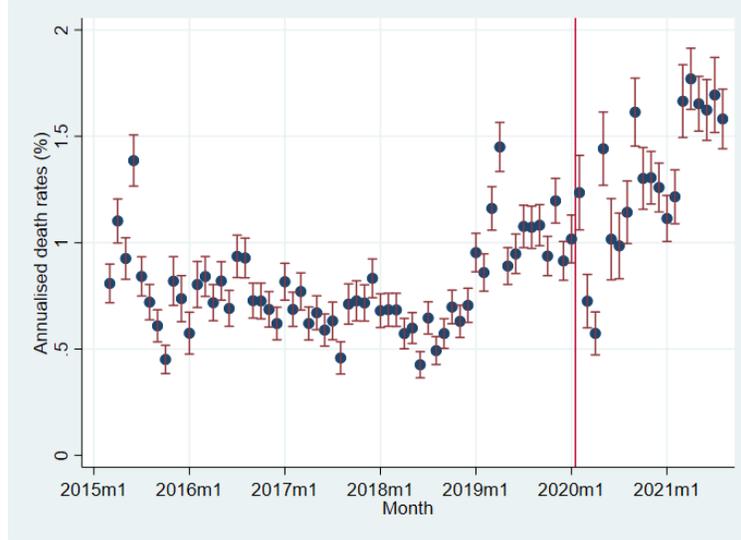
The annualized death rate calculated from the CPHS rises in 2019, a year before the pandemic (Figure A3). This is consistent with what is observed in the Global Burden of Disease data, though the magnitude is larger in the CPHS. This raises two questions. One, is there a pre-trend that begins in 2019? Perhaps the death rate during COVID is actually caused by something else started in 2019. Two, what years are the appropriate benchmarks against which excess deaths during the pandemic should be calculated?

We do not think that the jump in 2019 is a pre-trend unrelated to COVID that continues into 2020 for several reasons. First, there isn't a change in sample composition in 2019 leading to the rise in death rates. Of all households who respond in 2018, 98.5% households also respond in 2019. Also, of all households who respond in 2019, 99% households had also responded in 2018.

Second, as Figure 2 will show, the age-wise death rate in 2019 is significantly higher than the death rate that for 2015-2018 for both the elderly (age 60+) and for youth (0-17). By contrast, the age-wise death rate during the pandemic is significantly higher than during 2019 for only the elderly (60+). The jump in 2019 is not consistent with the age profile of COVID deaths, while the jump in 2020 is.

We use both 2019 and 2015-2019 as a baseline for the purposes of calculating excess deaths

Figure A3: Annualized death rates by month from January 2015 - June 2021.



Notes: The sample includes $t+1$ respondent and non-respondent households. Deaths reported in month t are allocated to month $t - 2$, for reasons explained later. The red line demarcates the start of the pandemic in February 2020.

during COVID. The argument for using 2019 is that the implied excess deaths is more in line with estimates based on CRS data after correcting for undercounting (e.g., Banaji and Gupta, 2021). The argument against using 2019 as a baseline is that it implies a baseline death rate of 1.07%, which is much higher than the death rate reported in the Global Burden of Disease. By contrast, the death rate implied by the 2015-2019 baseline (7.9%) is closer to the GBD baseline. Because the purpose of this paper is to estimate excess deaths and not the baseline death rate, we prefer employing the 2019 baseline, even though we report results from both baselines.

A natural question is whether the CPHS is to be believed given how high the baseline is in 2019. Our main answer is that, just because the baseline is high does not mean the change from 2019 to 2020-21 is incorrect. Indeed, our estimates of excess deaths from CPHS will be in line with the median estimate of excess deaths from the 10 states that have reported CRS data thus far.

B.4 Representativeness of CPHS

Dreze and Somanchi (2021) argues that CPHS undersamples the poor based on evidence that it yields both higher levels of literacy and faster improvement in literacy than government surveys. Dreze and Somanchi’s criticism is a problem for us if the poor have a different death rate during

COVID. The problem gets worse if the CPHS sample becomes less representative over time. It is possible that the poor have a lower death rate, as we show in Section 3.2.2; this would lower our estimate of excess deaths from COVID if the overall number overweights the rich. However, we do not believe that the problem gets worse over time as the gape between our control period (2019) and treatment period (2020-May 2021) is quite short. Moreover, experience from a serological survey conducted by one of us in Mumbai suggests that sampling the wealthy is more difficult than sampling the poor during lockdown (Malani et al., 2020).

Of course, the fact of difference between CPHS and government surveys of literacy is not dispositive of whether CPHS is biased since it is possible that the government surveys are the ones that are off. After all, the government has taken steps to suppress data (e.g., the 2017-18 consumption survey by the National Statistical Survey Office) that it finds unflattering. Moreover, government surveys are known to give different estimates of items like slums populations, with the differences driven by the policy aim of the survey.¹¹ Dreze and Somanchi suggest that CPHS is the one likely to be wrong because its frame samples more from the main streets of villages than from outskirts, where the poor tend to live. CMIE has responded that its sampling does get to outskirts and that the bias has not changed over time because that sampling frame is largely fixed and that its method for selection (of new households) has been constant (Vyas, 2021).

C Estimating infections

We estimate infections in three steps. First, we obtain data from a population-representative seroprevalence survey by the state of Tamil Nadu. The survey was conducted on 26,140 persons from 15 October - 30 November 2020. The sample included individuals aged 18 and above who provided informed consent. Details on the study and its estimate of seroprevalence are available from Malani et al. (2021). The survey was designed to be representative for urban and rural areas of districts and for demographic groups defined by age and gender.

Second, we extrapolate seroprevalence rates from this study to districts on 30 November 2020 based on the urban versus rural share of districts because urban share predicts more variability in seroprevalence than demographics and demographics do not vary much across districts.

Third, we take the curve that describes each district's new confirmed case over time (from

¹¹For example, public health officials in Mumbai in private conversations have noted that surveys of slum populations by the public health department tend to generate higher estimates of slum population because higher numbers in slums are more likely to generate large appropriations for the department.

www.covid19india.org) and scale it vertically up so that the total sum of scaled cases until 30 November equals the districts population times its estimated seroprevalence on that date. The resulting curve estimates the number of new individuals in a district infected with COVID on each day.

We had lots of choices for serological studies to use for our scaling. We could not use them all because they would yield inconsistent curves. Nor was there a clear way to combine them in a meta-analysis sense. We chose to use only the Tamil Nadu study for several reasons. First, it has a large sample size relative to other studies, e.g., the Karnataka study by Mohanan et al. (2021). It was one of a few state-wide or bigger serological surveys. Second, the study provided district-wise estimates unlike the national studies by the Indian Council for Medical Research Sharma (2021). While the Tamil Nadu study generates estimates of infection that are larger than the ICMR’s first 3 rounds of serological surveys, it generates estimates that are in line with ICMR’s fourth survey, which was conducted after wave 2. Third, we worked on the Tamil Nadu study, we knew the design and could vouch for its quality. We also had access to the data from that study so could better extrapolate from it to other districts.

D Comparing different pandemic definitions

Table A4 presents a matrix of estimates of the baseline and excess-death rate as we vary the definition of the baseline period (columns) and the pandemic period (rows).

Table A4: Robustness of excess-death rates to pandemic definition

	Nov 2018	Dec 2018	Jan 2019	Feb 2019	Mar 2019
Dec 2019	0.31	0.28	0.26	0.25	0.22
	0.99	1.02	1.05	1.06	1.08
Jan 2020	0.34	0.31	0.29	0.28	0.26
	0.99	1.02	1.04	1.05	1.07
Feb 2020	0.36	0.33	0.31	0.30	0.28
	0.99	1.02	1.04	1.05	1.06
Mar 2020	0.35	0.33	0.31	0.30	0.28
	1.00	1.02	1.04	1.05	1.07
Apr 2020	0.38	0.36	0.34	0.33	0.32
	0.99	1.01	1.03	1.04	1.06

Notes. This table presents a matrix of estimates of the baseline and excess-death rate as we vary the definition of the baseline period (columns) and the pandemic period (rows). The columns refer to the start date of the baseline. The rows refer to the start date of the pandemic. Each cell contains an estimate of the annualized excess-death rate (top) and baseline death rate (bottom) based on (1).

E Second solution to the timing-of-death problem

The timing-of-death problem is illustrated with the example in Appendix Table A5. Suppose the true death rates over a 10 month period are 7 per 1000 for each month except month 4, where it jumps to 9 (row 1). Moreover, assume one quarter of households are surveyed each month and all households respond to surveys, which are 4 months apart. Because only a quarter of households are interviewed each month, those extra 2 deaths will, statistically, be distributed over 4 months after the death (row 2). The problem we face is how to back out the jump to 9 in row 1.

Our preferred solution from the main text, which is illustrated in row 3 of Table A5, is akin to treating the reported number at t as a moving average of the true death rate for $t - (k/2)$ (row 4).

Here we explore a second solution: to estimate the death rate by asking the question, how much would the true death rate have to have changed for the observed death rate to have changed as much as it did since the last month. The formula that provides the answer is in row 5.

The advantage of the second solution is that it can, in some cases, back out the true death rate. But there are two problems. First, death rates need to be stable for a period otherwise one cannot solve the formula because it uses prior value of estimated rates to measure current rates. Second, if observed deaths have some error unrelated to timing, this solution magnifies those errors. The solution recognizes that the actual changes have to be larger than observed changes because the survey process smooths out changes over few months (row 2). But if there are errors in observed data, then the errors are also magnified. This can increase variability of results from solution 2, which we will demonstrate. Because we think there could be errors in CPHS rosters, our preferred measure is the first one.

We implement our second solution to the timing-of-death problem with a regression of the form

$$d_{i,t,t-k} = \sum_{k=4,8,\dots} \delta_k \cdot \mathbb{I}(k) + \sum_{m \in \text{pandemic}} \beta_m \cdot \mathbb{I}(t-k \leq m \leq t) + \epsilon_{ist} \quad (2)$$

where $\mathbb{I}(k)$ is an indicator for whether the gap between the current survey and the one to which she last responded is k months, and $\mathbb{I}(t-k \leq m \leq t)$ is an indicator whether the intervening period between surveys was during the pandemic. Standard errors are clustered at the village/ward \times month level to account for correlation in reporting of deaths within a locality.

The coefficient δ_k estimates a pre-COVID death rate for observations that are k months apart and our parameter of interest β_m captures the increment in true death rate implied by the increment

Table A5: Timing of CPHS observation of deaths and correction for that timing

Formula		Annualized death rate (deaths/1000) in each month											
		1	2	3	4	5	6	7	8	9	10	11	12
1. Truth (d)		7	7	7	7	9	7	7	7	7	7	7	7
2. Observed (y)	$y_t = (d_t + d_{t-1} + d_{t-2} + d_{t-3})/4$	7	7	7	7	7.5	7.5	7.5	7.5	7	7	7	7
3. Solution 1 (z)	$z_t = y_{t-2}$	7	7	7.5	7.5	7.5	7.5	7	7	7	7	7	7
4. Moving average of d	$m_t = (d_{t-1} + \dots + y_{t+2})/4$	7	7	7.5	7.5	7.5	7.5	7	7	7	7	7	7
5. Solution 2	$r_t = (4y_t) - (r_{t-1} + r_{t-2} + r_{t-3})$	7	7	7	7	9	7	7	7	7	7	7	7

in observed death rate in each month during COVID. The coefficient γ_i nets out seasonality in deaths.

Our second solution to the timing-of-death problem produces implausibly variable and/or high estimates of excess deaths from the pandemic.

Excess variability is illustrated in Figure A4 Panel A, which plots monthly excess-death rates (relative to 2019 rates) based on (a) estimating (2), but with month fixed effects instead of a pandemic period indicator, and (b) using, importantly, observations only from households that answer consecutive surveys. This solution produces a strong saw-tooth pattern with spikes every 4 months in the excess-death rate. The reason for this pattern is that the second solution is premised on the idea that a jump in month t is reallocated evenly across months t to $t+3$. The solution corrects that by reallocating deaths from months 2-4 back to month 1. But a jump on observed error may reflect a true increase in death rates, or an error in the roster. If the jump were a true jump in death rates, that correction would be correct. But if it a positive error in month t , then months $t+1$ to $t+3$ are inappropriately suppressed. The suppression ends disappears in month $t+4$ so month $t+4$ is higher than $t+3$. If there happens to be another positive shock in $t+4$, the pattern repeats. We believe that is what was happening in Figure A4.

Including observations from households that may not answer consecutive surveys mitigates the saw-tooth pattern (Figure A4 Panel B). Consecutive surveys are 4 months apart and, so, solution 2 reallocates positive jumps in deaths only over 4 months. When that is relaxed, some of the jump is allocated over 4 months, some over 8 months, and so on. This dampens the 4-month cycles our second solution generates from responders to consecutive surveys.

Estimates of excess deaths during the pandemic are higher when we use our second solution than in our first solution. Table A6 reports our estimate of excess deaths using solution 2 for various definitions of the pandemic period. If we only use observations from households that answer consecutive observations, our estimate of the annualized pandemic period excess-death rate

is 1.763 and significant if we use our preferred pandemic start date (February 2020). It falls to 0.497 and insignificant if we use all observations. The latter is partly due to the fact that deaths in households that do not answer the last round of 2020 or the first round of 2021 are included but we have not captured deaths in those households yet; they may report these deaths in future rounds. As with our first solution, the excess-death rate under the second solution falls as we move up the start date and rises as we move it back.

Table A6: excess-death rates under the second timing solution for different pandemic definitions

Panel A: Consecutive observations

Pandemic start month	Excess death rates (Annualised rate (%))	Standard error
Dec 2019	2.229	(0.572)
Jan 2020	2.145	(0.594)
Feb 2020	2.712	(0.578)
Mar 2020	2.871	(0.606)
Apr 2020	2.952	(0.639)
May 2020	3.291	(0.675)

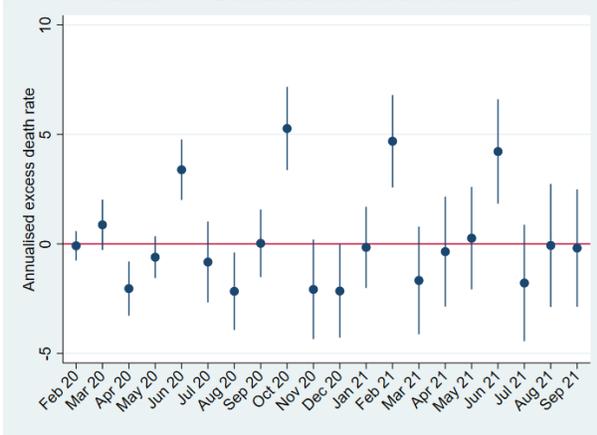
Panel B: All observations

Pandemic start month	Excess death rates (Annualised rate (%))	Standard error
Dec 2019	1.003	(0.501)
Jan 2020	0.885	(0.521)
Feb 2020	1.154	(0.495)
Mar 2020	1.456	(0.519)
Apr 2020	1.662	(0.547)
May 2020	2.125	(0.577)

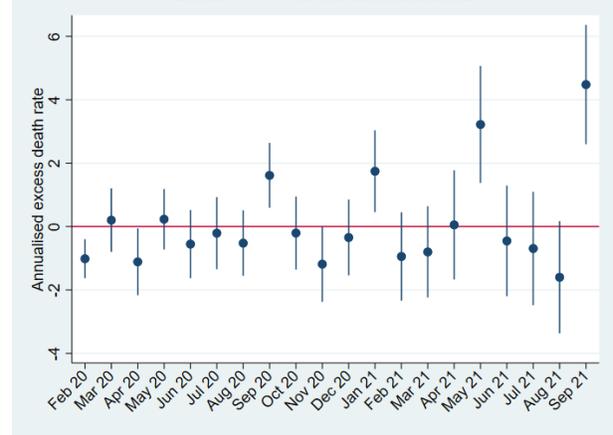
Notes. This table reports the excess-death rates computed using the specification in equation 2. The different rows represent different start dates for the pandemic period. The excess-death rates are the mean of the monthly death rates multiplied by 12. The table on the left only uses observations from households that answer consecutive surveys. The table on the right uses all observations.

Figure A4: Monthly death rates from second solution

Panel A: Consecutive observations



Panel B: All observations



Notes. This figure plots monthly death rates for each month. The coefficients from regression 2 are plotted here. The indicator for month t is 1 for an observation if the month t is between the month in which the individual was surveyed (including month of survey) and the month in which the individual is next surveyed in (excluding the month of survey). The regression is weighted using the individual's weights. Results in Panel A use only responses from households that answer consecutive surveys. Results in Panel B use responses from all households.

F Age-wise deaths

Table A7 reports estimates of a regression based on (1), but on samples in different age bins. It shows that the excess-death rate during the pandemic is larger and significant in older ages. Moreover, this age skew was more pronounced in the second wave.

Table A7: Excess-death rate during the pandemic and its waves in different age groups

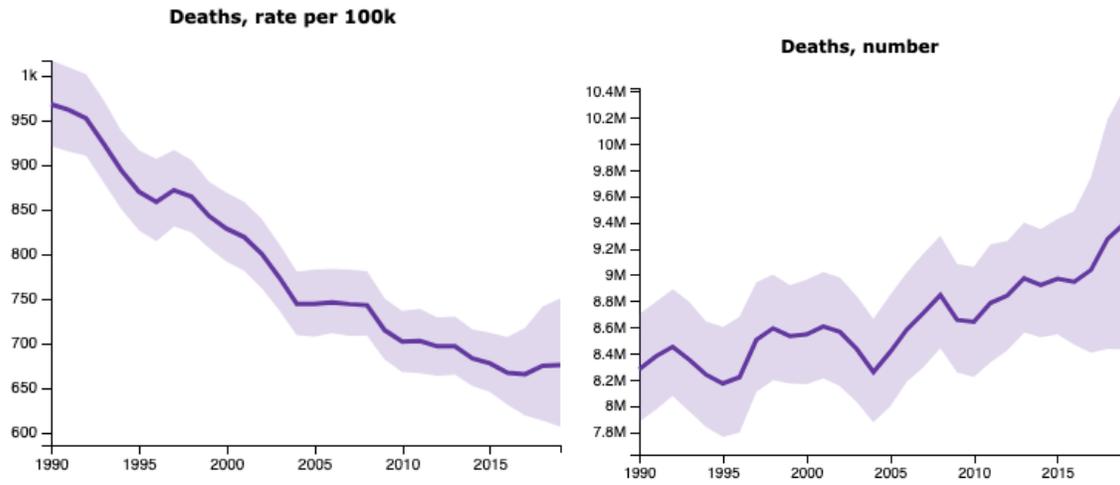
<i>Panel A: Pandemic overall</i>							
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
	0-17	18-29	30-39	40-49	50-59	60-69	70+
During Pandemic	-0.0526 (0.0583)	-0.0389 (0.0414)	-0.0313 (0.0512)	0.0286 (0.0496)	0.245* (0.118)	1.865*** (0.354)	6.164*** (0.989)
2019 mean	0.355*** (0.0423)	0.318*** (0.0324)	0.319*** (0.0401)	0.405*** (0.0344)	1.229*** (0.0769)	4.175*** (0.204)	12.98*** (0.572)
N	618580	606601	370133	484492	369772	175608	82468

<i>Panel B: By waves</i>							
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
	0-17	18-29	30-39	40-49	50-59	60-69	70+
Wave 1	-0.0176 (0.0708)	0.00249 (0.0489)	-0.0467 (0.0556)	0.0302 (0.0584)	0.207 (0.135)	1.277** (0.441)	2.162** (0.824)
Wave 2	-0.101 (0.0616)	-0.0944* (0.0448)	-0.00975 (0.0695)	0.0265 (0.0660)	0.295 (0.153)	2.609*** (0.459)	10.93*** (0.897)
2019 mean	0.355*** (0.0423)	0.318*** (0.0324)	0.319*** (0.0401)	0.405*** (0.0344)	1.229*** (0.0769)	4.175*** (0.204)	12.98*** (0.393)
N	618580	606601	370133	484492	369772	175608	84692

Notes. The sample includes households that respond in consecutive rounds. Observations on individuals are weighted to be nationally representative even with non-response. The regression model in Panel A is (1); the model in Panel B replaces the pandemic indicator in (1) with wave indicators. Each column reports estimates from a different regression. Regressions for each category are run by restricting the sample to those in that age category alone. Standard errors clustered at the village/ward \times month level are reported in parentheses. */**/** indicates $p < 0.05/0.01/0.001$.

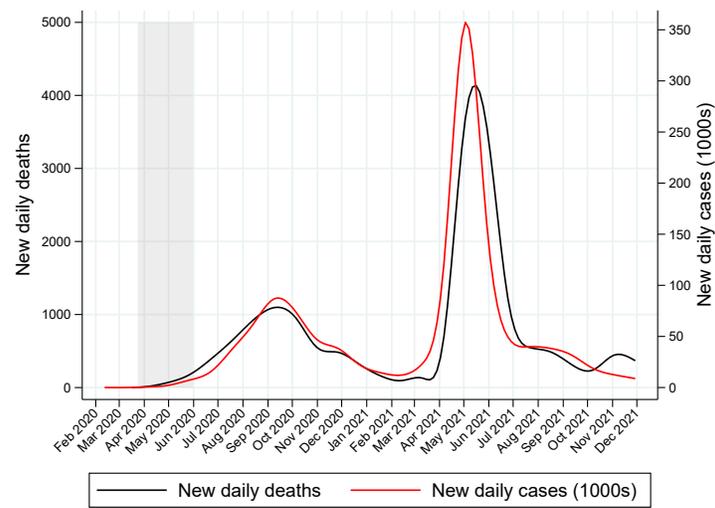
G Appendix exhibits

Figure A5: Death rate and total deaths in India over time.



Source: ghdx.healthdata.org/gbd-results-tool.

Figure A6: Confirmed COVID cases and deaths over time.



Source: www.covid19india.org. We thank Bartek Woda for making this figure.

Table A8: COVID death rate prior to the pandemic and the excess-death rate during the pandemic

	Annualised death rates (%)				
	(1)	(2)	(3)	(4)	(5)
During Pandemic	0.543*** (0.0432)	0.292*** (0.0504)	0.275*** (0.0506)	0.279*** (0.0506)	
Baseline	0.787*** (0.0128)	1.038*** (0.0304)	2.416*** (0.387)	1.953*** (0.373)	1.038*** (0.0304)
Wave 1					0.153** (0.0559)
Wave 2					0.478*** (0.0739)
Controls	None	None	HR FE	District FE	None
Data used	2015-21	2019-21	2019-21	2019-21	2019-21
N	7991330	2969850	2969850	2969850	2969850

Notes. The sample includes households that respond in consecutive rounds. Observations on individuals are weighted to be nationally representative even with non-response. The regression model is (1). HR FE means fixed effects for homogeneous region, a cluster of similar districts within a state; district FE means district fixed effects. Standard errors clustered at the village/ward \times month level are reported in parentheses. In the column 2, if we instead cluster standard errors at the homogeneous region \times month level, the standard error is 0.073 and the effect remains significant. */**/** indicates $p < 0.05/0.01/0.001$.

Table A9: Annualised death rates by gender and urban status

<i>Panel A: Gender</i>			<i>Panel B: Urban/rural</i>		
	Annualised death rates(%)			Annualised death rates(%)	
	(1)	(2)		(1)	(2)
Pandemic	0.285*** (0.0578)		Pandemic	0.239*** (0.0587)	
Wave 1		0.142* (0.0638)	Wave 1		0.176* (0.0730)
Wave 2		0.477*** (0.0817)	Wave 2		0.326*** (0.0719)
Pandemic × Female	0.0148 (0.0635)		Pandemic × Urban	0.160 (0.110)	
Wave 1 × Female		0.0235 (0.0798)	Wave 1 × Urban		-0.0691 (0.109)
Wave 2 × Female		0.00239 (0.0869)	Wave 2 × Urban		0.435* (0.173)
Female	-0.000111 (0.000128)	-0.000111 (0.000128)	Urban	-0.000105 (0.000215)	-0.000105 (0.000215)
2019 mean	1.054*** (0.0356)	1.054*** (0.0356)	2019 mean	1.048*** (0.0369)	1.048*** (0.0369)
<i>N</i>	2969850	2969850	<i>N</i>	2969850	2969850

Notes. Estimates are from a regression model based on 1, with the addition of a gender (urban) indicator and gender (urban) indicator interacted with the pandemic or wave indicator. Sample includes only consecutive observations and excludes emigrants. Standard errors clustered at the village/ward × month level are reported in parentheses. $p < 0.05/0.01/0.001$.