

Drivers of COVID-19 deaths in the United States: A two-stage modeling approach

Christopher F Baum (Boston College, DIW Berlin & CESIS)

Andrés García-Suaza (Universidad del Rosario, Bogotá)

Jesús Otero (Universidad del Rosario, Bogotá)

Miguel Henry (Greylock McKinnon Associates)

2023 UK Stata Conference, London, September 2023

Introduction

We present a two-stage econometric modeling approach to examine a number of socioeconomic, demographic, health, epidemiological, climate and political drivers affecting the spread of COVID-19 across six pandemic waves for counties in the United States.

Our empirical strategy exploits the availability of two years of daily data: March 15, 2020 through March 19, 2022 on the number of confirmed deaths and cases of COVID-19 in 3,014 U.S. counties of the 48 contiguous states and the District of Columbia. We also make use of daily county-level vaccination rate data for the period in which vaccinations were available.

Introduction

We present a two-stage econometric modeling approach to examine a number of socioeconomic, demographic, health, epidemiological, climate and political drivers affecting the spread of COVID-19 across six pandemic waves for counties in the United States.

Our empirical strategy exploits the availability of two years of daily data: March 15, 2020 through March 19, 2022 on the number of confirmed deaths and cases of COVID-19 in 3,014 U.S. counties of the 48 contiguous states and the District of Columbia. We also make use of daily county-level vaccination rate data for the period in which vaccinations were available.

In the first stage of the analysis, we use a daily-frequency panel data set on COVID-19 cases and deaths to fit mixed models of cases and deaths against lagged confirmed cases and lagged COVID-19 vaccinations for each county.

As the resulting intercept and slope coefficients are county-specific, they relax the homogeneity assumption that is implicit when the analysis is performed using geographically aggregated cross-section units.

In the first stage of the analysis, we use a daily-frequency panel data set on COVID-19 cases and deaths to fit mixed models of cases and deaths against lagged confirmed cases and lagged COVID-19 vaccinations for each county.

As the resulting intercept and slope coefficients are county-specific, they relax the homogeneity assumption that is implicit when the analysis is performed using geographically aggregated cross-section units.

In the second stage of the analysis, we assume that the county-level slope coefficient point estimates are a function of factors that are taken as fixed over the course of the pandemic. As these are generated data, we take their precision into account.

To guide the choice of regressors in the second stage, we employ the novel one-covariate-at-a-time variable selection OCMT algorithm proposed by Chudik, Kapetanios, and Pesaran (2018).

To contrast the importance of factors over the six pandemic waves, we employ an unorthodox approach based on the seemingly unrelated regression (SUR) model.

In the second stage of the analysis, we assume that the county-level slope coefficient point estimates are a function of factors that are taken as fixed over the course of the pandemic. As these are generated data, we take their precision into account.

To guide the choice of regressors in the second stage, we employ the novel one-covariate-at-a-time variable selection OCMT algorithm proposed by Chudik, Kapetanios, and Pesaran (2018).

To contrast the importance of factors over the six pandemic waves, we employ an unorthodox approach based on the seemingly unrelated regression (SUR) model.

In the second stage of the analysis, we assume that the county-level slope coefficient point estimates are a function of factors that are taken as fixed over the course of the pandemic. As these are generated data, we take their precision into account.

To guide the choice of regressors in the second stage, we employ the novel one-covariate-at-a-time variable selection OCMT algorithm proposed by Chudik, Kapetanios, and Pesaran (2018).

To contrast the importance of factors over the six pandemic waves, we employ an unorthodox approach based on the seemingly unrelated regression (SUR) model.

Related literature

There has been an explosion of studies of the COVID-19 pandemic from its advent to the present day. Many of those studies have considered the geography of the pandemic's spread, for the U.S. and other countries. Many have also focused on demographic and socioeconomic factors to the degree that they vary across geography. Most of the cited literature addresses the early phases of the pandemic. In our study, we analyze the pandemic's evolution over a longer period.

For brevity, I do not describe the studies here., but they are included in the references.

Motivation for analysis by waves

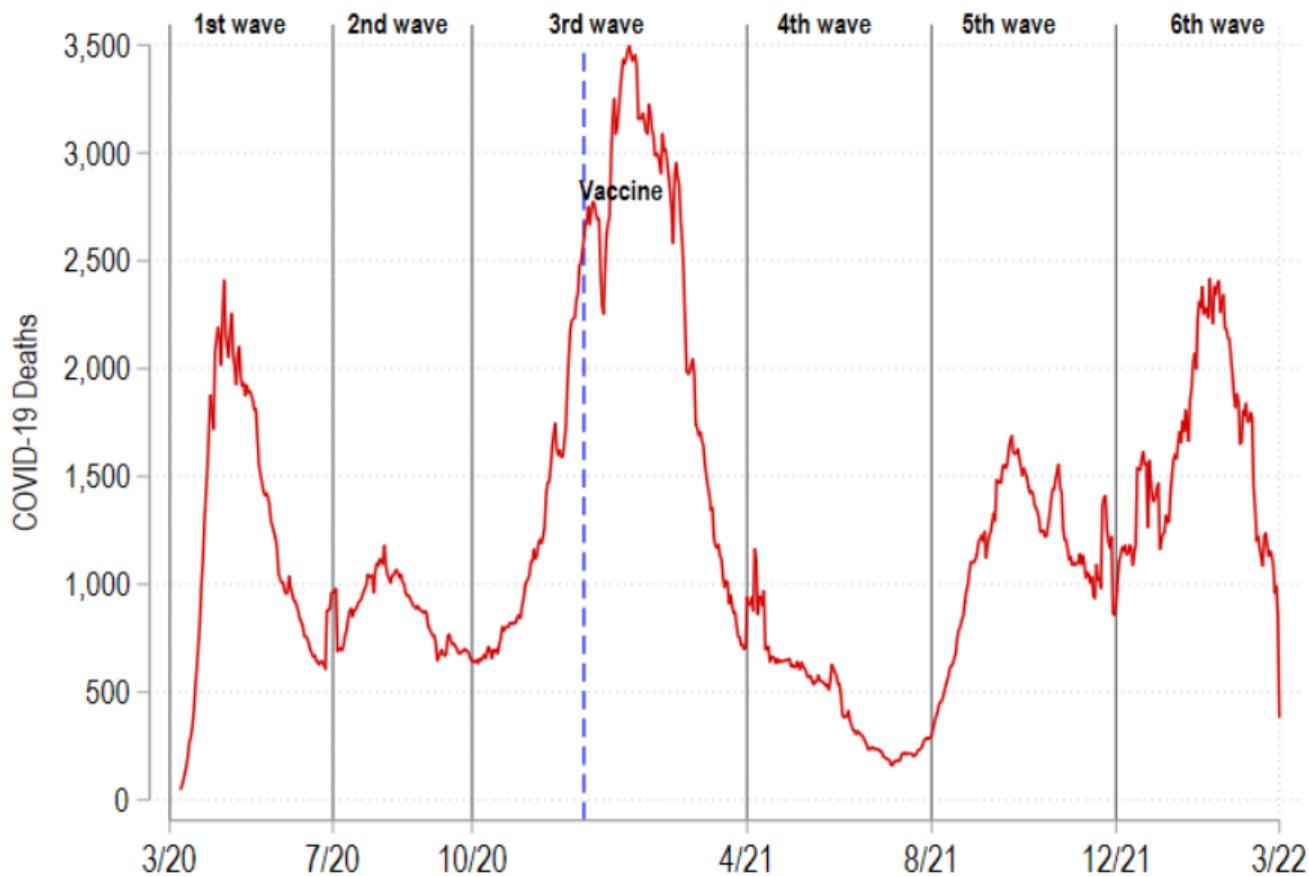
To model the evolution of the pandemic in the U.S., we recognize that its severity has varied considerably, given the mutating dominant variants, the introduction of widespread vaccinations, and improvements in treatment of the disease such as immediate treatment of an infection with Paxlovid. The latter factor is particularly important as it has reduced the likelihood of mortality for those infected in most segments of the population.

Unlike many of the earlier studies, our analysis includes two full years of daily data, and must consider the temporal stability of an estimated model in the context of these time-varying factors affecting the severity of the virus and the ability to prevent infections and treat them.

A single model is not adequate to capture these variations over the past two years in the U.S. We adopt the nomenclature used by the Pew Research Center in Jones (2022), which identifies six distinct waves. The following figure shows the trajectory of deaths attributed to COVID-19. Each wave is identified by its starting month. The dashed vertical line in mid-December 2020 denotes the introduction of vaccinations.

Unlike many of the earlier studies, our analysis includes two full years of daily data, and must consider the temporal stability of an estimated model in the context of these time-varying factors affecting the severity of the virus and the ability to prevent infections and treat them.

A single model is not adequate to capture these variations over the past two years in the U.S. We adopt the nomenclature used by the Pew Research Center in Jones (2022), which identifies six distinct waves. The following figure shows the trajectory of deaths attributed to COVID-19. Each wave is identified by its starting month. The dashed vertical line in mid-December 2020 denotes the introduction of vaccinations.



The six waves include daily county-level data on confirmed cases and deaths for the periods:

- 1: 15 March 2020 - 30 June 2020
- 2: 1 July 2020 - 30 September 2020
- 3: 1 October 2020 - 31 March 2021
- 4: 1 April 2021 - 31 July 2021
- 5: 1 August 2021 - 30 November 2021
- 6: 1 December 2021 - 19 March 2022

Wave 3 partially captures the rollout of vaccines in mid-December 2020, while wave 5 reflects the surge of the delta variant, and wave 6 displays the dominance of the omicron variant.

The six waves include daily county-level data on confirmed cases and deaths for the periods:

- 1: 15 March 2020 - 30 June 2020
- 2: 1 July 2020 - 30 September 2020
- 3: 1 October 2020 - 31 March 2021
- 4: 1 April 2021 - 31 July 2021
- 5: 1 August 2021 - 30 November 2021
- 6: 1 December 2021 - 19 March 2022

Wave 3 partially captures the rollout of vaccines in mid-December 2020, while wave 5 reflects the surge of the delta variant, and wave 6 displays the dominance of the omicron variant.

The six waves include daily county-level data on confirmed cases and deaths for the periods:

- 1: 15 March 2020 - 30 June 2020
- 2: 1 July 2020 - 30 September 2020
- 3: 1 October 2020 - 31 March 2021
- 4: 1 April 2021 - 31 July 2021
- 5: 1 August 2021 - 30 November 2021
- 6: 1 December 2021 - 19 March 2022

Wave 3 partially captures the rollout of vaccines in mid-December 2020, while wave 5 reflects the surge of the delta variant, and wave 6 displays the dominance of the omicron variant.

The six waves include daily county-level data on confirmed cases and deaths for the periods:

- 1: 15 March 2020 - 30 June 2020
- 2: 1 July 2020 - 30 September 2020
- 3: 1 October 2020 - 31 March 2021
- 4: 1 April 2021 - 31 July 2021
- 5: 1 August 2021 - 30 November 2021
- 6: 1 December 2021 - 19 March 2022

Wave 3 partially captures the rollout of vaccines in mid-December 2020, while wave 5 reflects the surge of the delta variant, and wave 6 displays the dominance of the omicron variant.

The six waves include daily county-level data on confirmed cases and deaths for the periods:

- 1: 15 March 2020 - 30 June 2020
- 2: 1 July 2020 - 30 September 2020
- 3: 1 October 2020 - 31 March 2021
- 4: 1 April 2021 - 31 July 2021
- 5: 1 August 2021 - 30 November 2021
- 6: 1 December 2021 - 19 March 2022

Wave 3 partially captures the rollout of vaccines in mid-December 2020, while wave 5 reflects the surge of the delta variant, and wave 6 displays the dominance of the omicron variant.

The six waves include daily county-level data on confirmed cases and deaths for the periods:

- 1: 15 March 2020 - 30 June 2020
- 2: 1 July 2020 - 30 September 2020
- 3: 1 October 2020 - 31 March 2021
- 4: 1 April 2021 - 31 July 2021
- 5: 1 August 2021 - 30 November 2021
- 6: 1 December 2021 - 19 March 2022

Wave 3 partially captures the rollout of vaccines in mid-December 2020, while wave 5 reflects the surge of the delta variant, and wave 6 displays the dominance of the omicron variant.

The six waves include daily county-level data on confirmed cases and deaths for the periods:

- 1: 15 March 2020 - 30 June 2020
- 2: 1 July 2020 - 30 September 2020
- 3: 1 October 2020 - 31 March 2021
- 4: 1 April 2021 - 31 July 2021
- 5: 1 August 2021 - 30 November 2021
- 6: 1 December 2021 - 19 March 2022

Wave 3 partially captures the rollout of vaccines in mid-December 2020, while wave 5 reflects the surge of the delta variant, and wave 6 displays the dominance of the omicron variant.

The six waves include daily county-level data on confirmed cases and deaths for the periods:

- 1: 15 March 2020 - 30 June 2020
- 2: 1 July 2020 - 30 September 2020
- 3: 1 October 2020 - 31 March 2021
- 4: 1 April 2021 - 31 July 2021
- 5: 1 August 2021 - 30 November 2021
- 6: 1 December 2021 - 19 March 2022

Wave 3 partially captures the rollout of vaccines in mid-December 2020, while wave 5 reflects the surge of the delta variant, and wave 6 displays the dominance of the omicron variant.

The cumulative cases and deaths and the vaccination rates for each of the six waves are presented in the following table. We also computed these measures for two subsets of counties: those in the 4th quartile of population density, labeled High, and those in the other three quartiles, labeled Low. The impact of population density on both cases and deaths is meaningful, particularly in the earlier waves.

Table: Cumulative cases, deaths, and vaccination rates (N = 2,215,290)

Wave starting:		3/20	7/20	10/20	4/21	8/21	12/21
Cases/100K	Total	522.71	1962.47	9394.12	10595.04	15728.08	23992.87
	Low	479.73	1959.10	9573.60	10731.38	16117.79	24172.88
	High	651.76	1972.59	8855.23	10185.63	14557.91	23452.38
Deaths/100K	Total	17.65	43.95	189.10	211.17	284.94	357.31
	Low	13.83	41.52	199.08	222.26	303.84	379.28
	High	29.12	51.26	159.13	177.85	228.16	291.35
Vaccinations	Total	0.00	0.00	13.66	32.65	45.54	50.90
	Low	0.00	0.00	13.71	30.89	43.42	48.47
	High	0.00	0.00	13.49	37.93	51.91	58.20

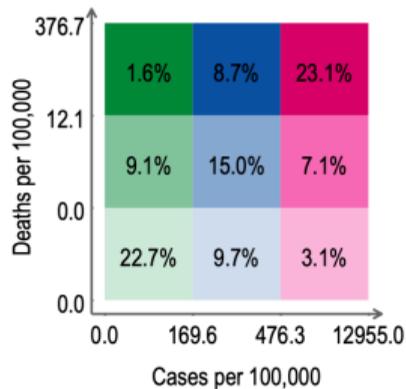
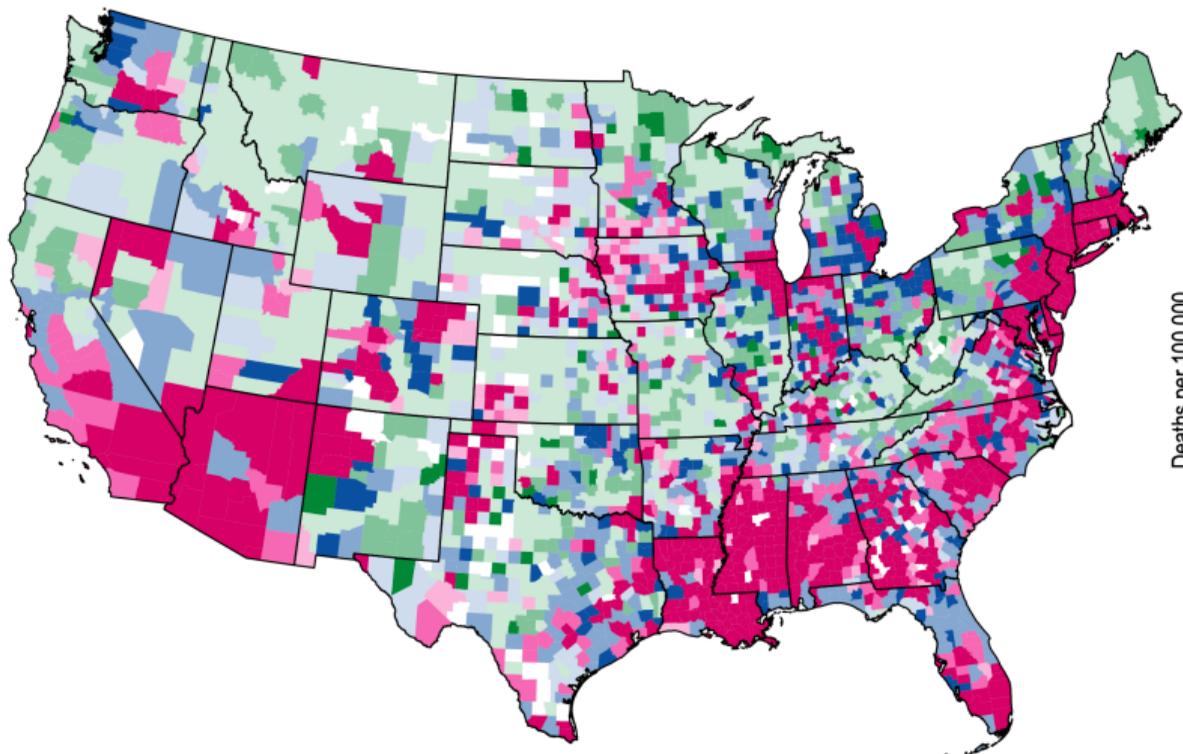
To visualize the variations in the COVID-19 cumulative cases and deaths, we consider how these variables are correlated across US counties. That visualization can be implemented by Asjad Naqvi's innovative `bimap` package (Naqvi (2022)), available from the SSC Archive and documented in his Medium guide for Bi-variate maps.

We present bivariate maps of these two variables' averages for each wave over the continental United States. The deep red color identifies the counties which are in the upper tercile of both case rates and death rates, while light green in the lower left identifies case rates and death rates in the first tercile.

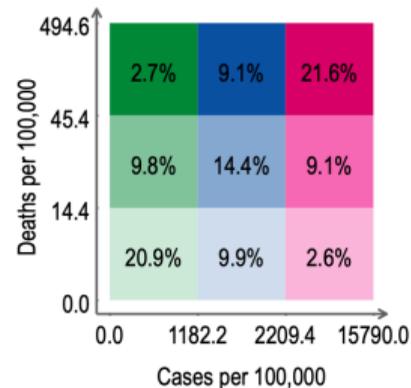
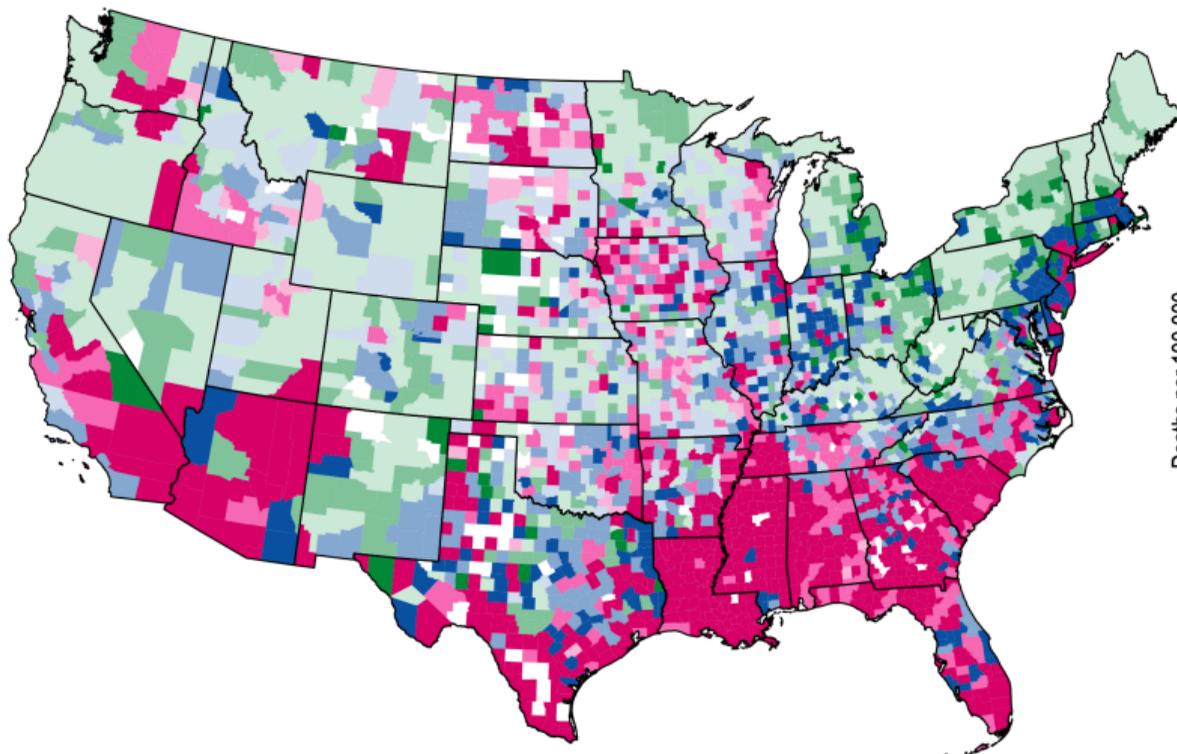
To visualize the variations in the COVID-19 cumulative cases and deaths, we consider how these variables are correlated across US counties. That visualization can be implemented by Asjad Naqvi's innovative `bimap` package (Naqvi (2022)), available from the SSC Archive and documented in his Medium guide for Bi-variate maps.

We present bivariate maps of these two variables' averages for each wave over the continental United States. The deep red color identifies the counties which are in the upper tercile of both case rates and death rates, while light green in the lower left identifies case rates and death rates in the first tercile.

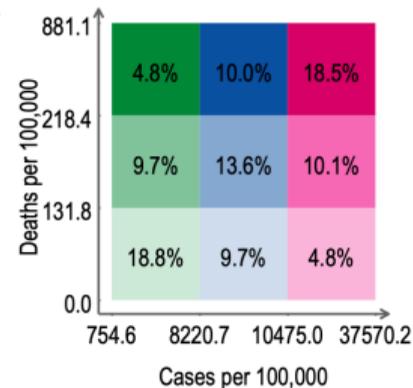
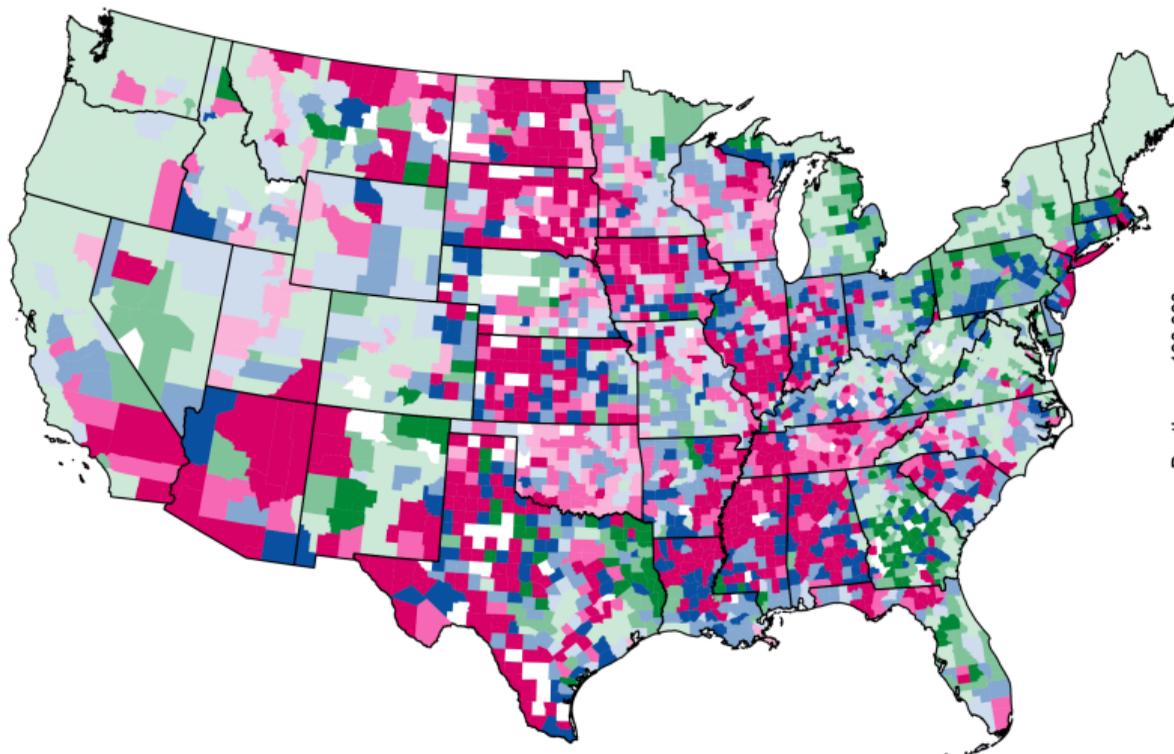
Wave 1: 15 March–30 June 2020



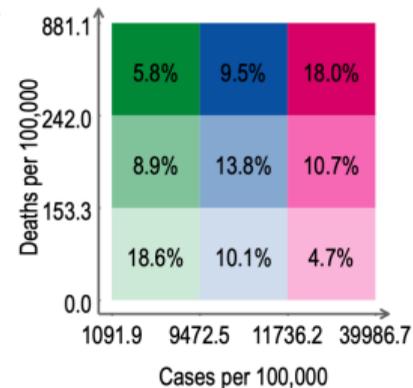
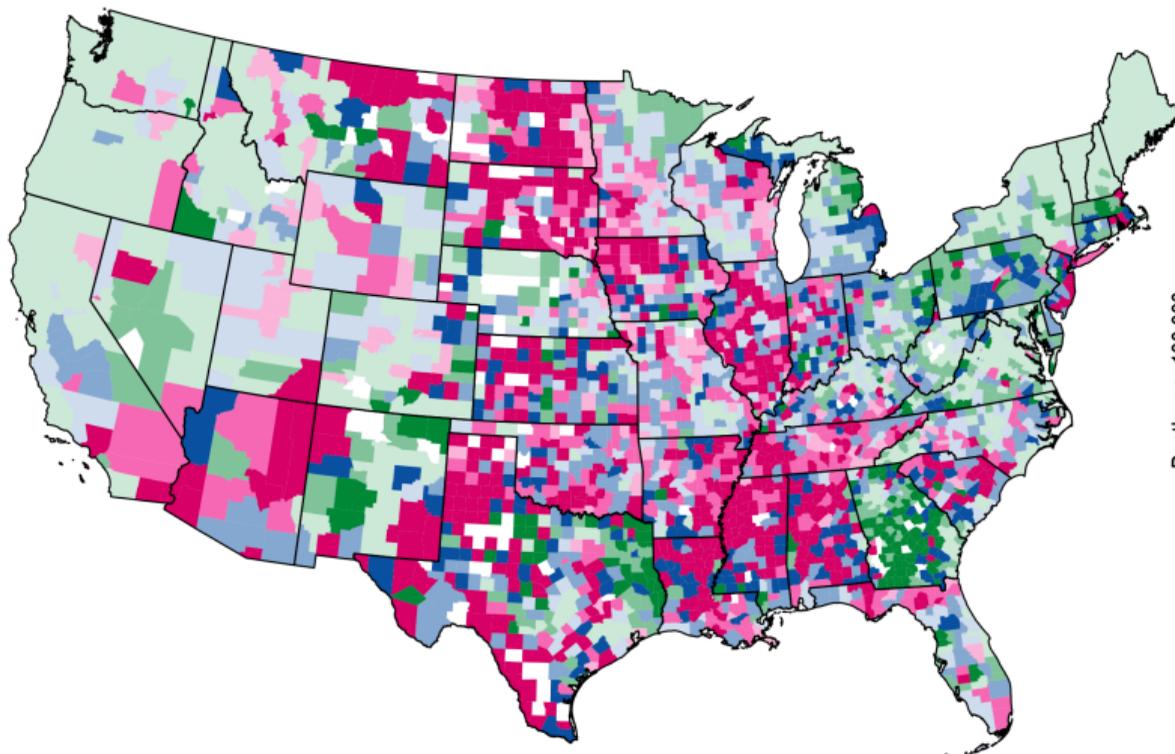
Wave 2: 1 July–30 September 2020



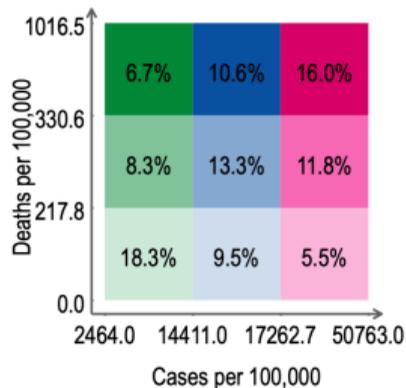
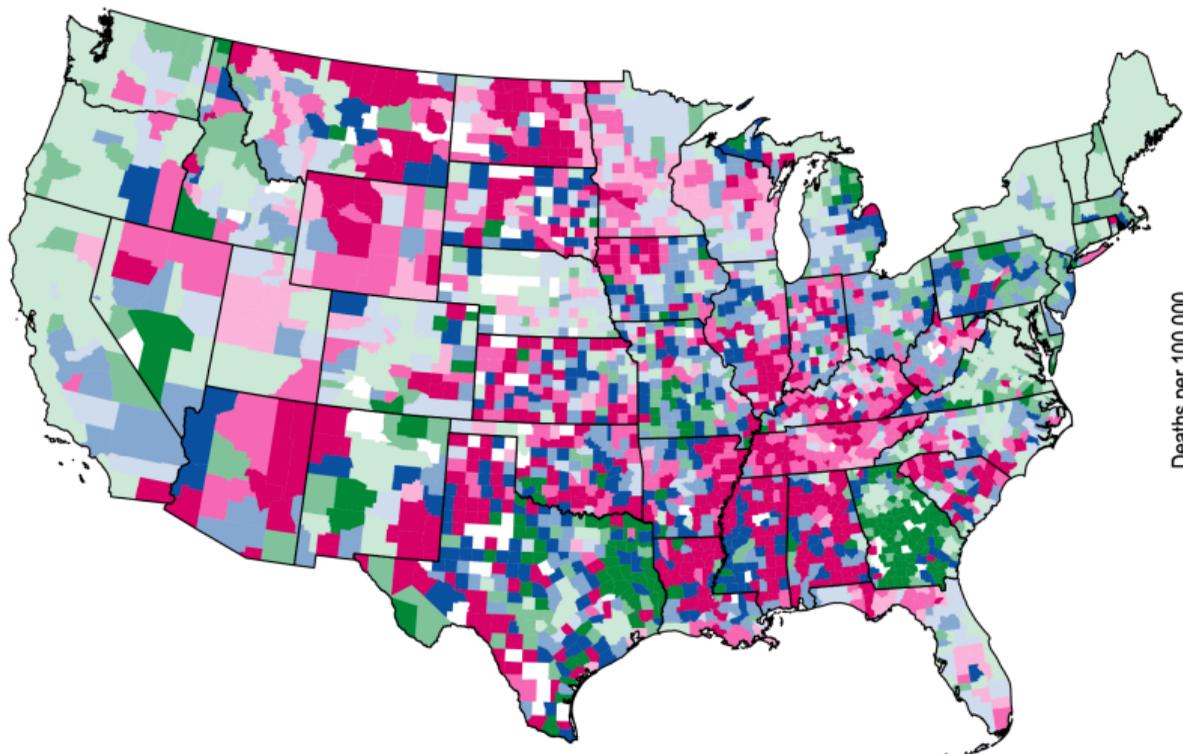
Wave 3: 1 October 2020–31 March 2021



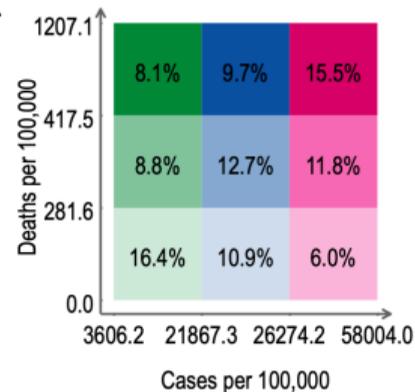
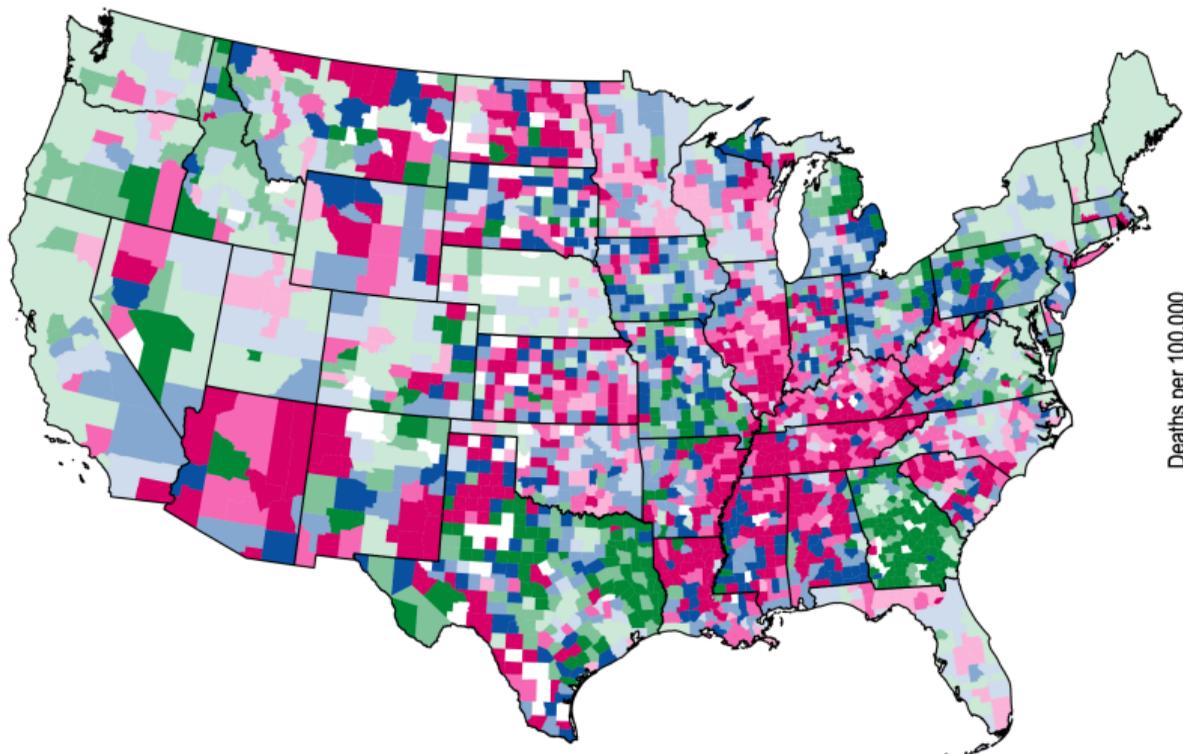
Wave 4: 1 April–31 July 2021



Wave 5: 1 August–30 November 2021



Wave 6: 1 December 2021–19 March 2022



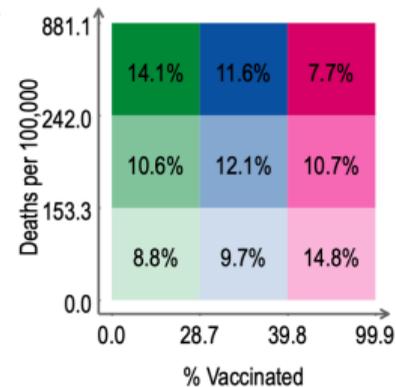
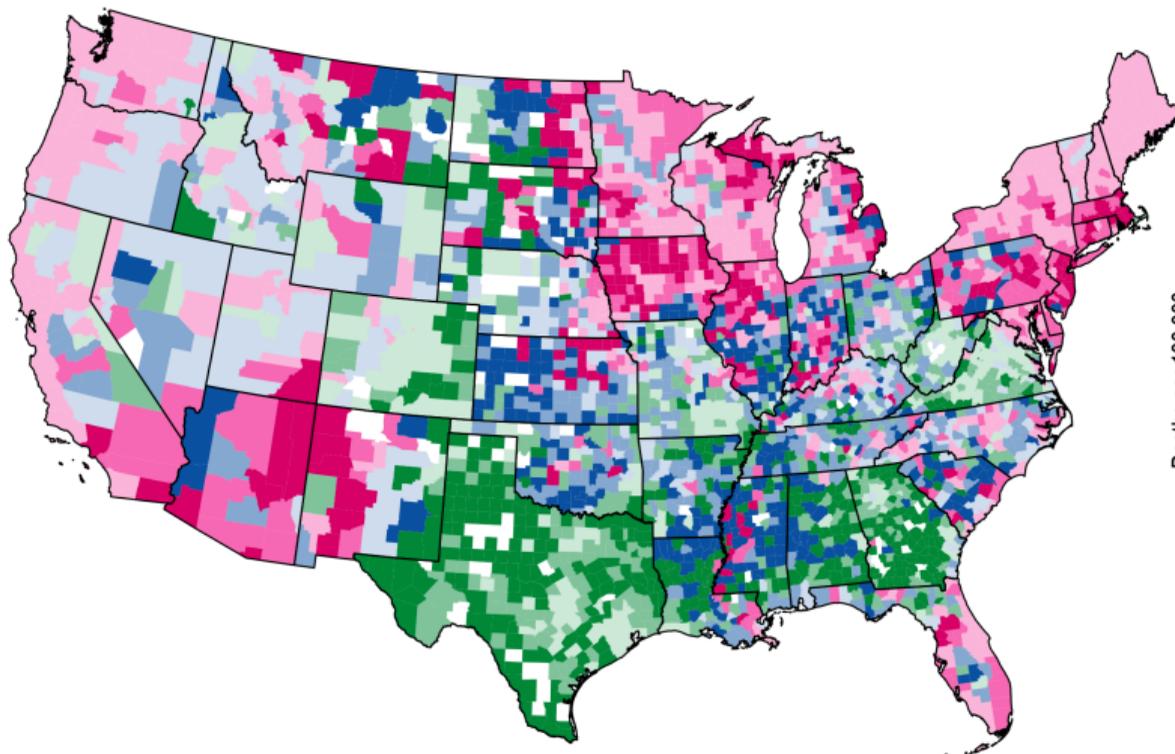
The following maps illustrate the bivariate relationship between vaccination rates and cumulative deaths in waves 4, 5 and 6. Vaccines first became available during wave 3, only being widespread over the first several months of 2021.

In the maps' legend, the rightmost categories refer to the third tercile of vaccination rates. The dark green counties are those with low vaccination rates and the third tercile of cumulative deaths. The geography of those most seriously affected by the course of the virus is quite evident.

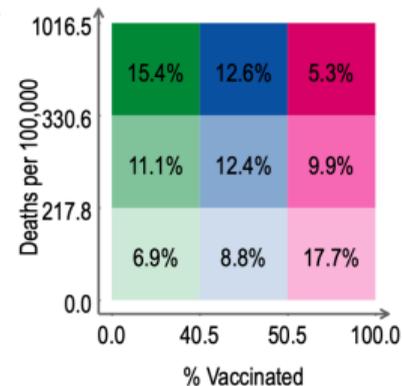
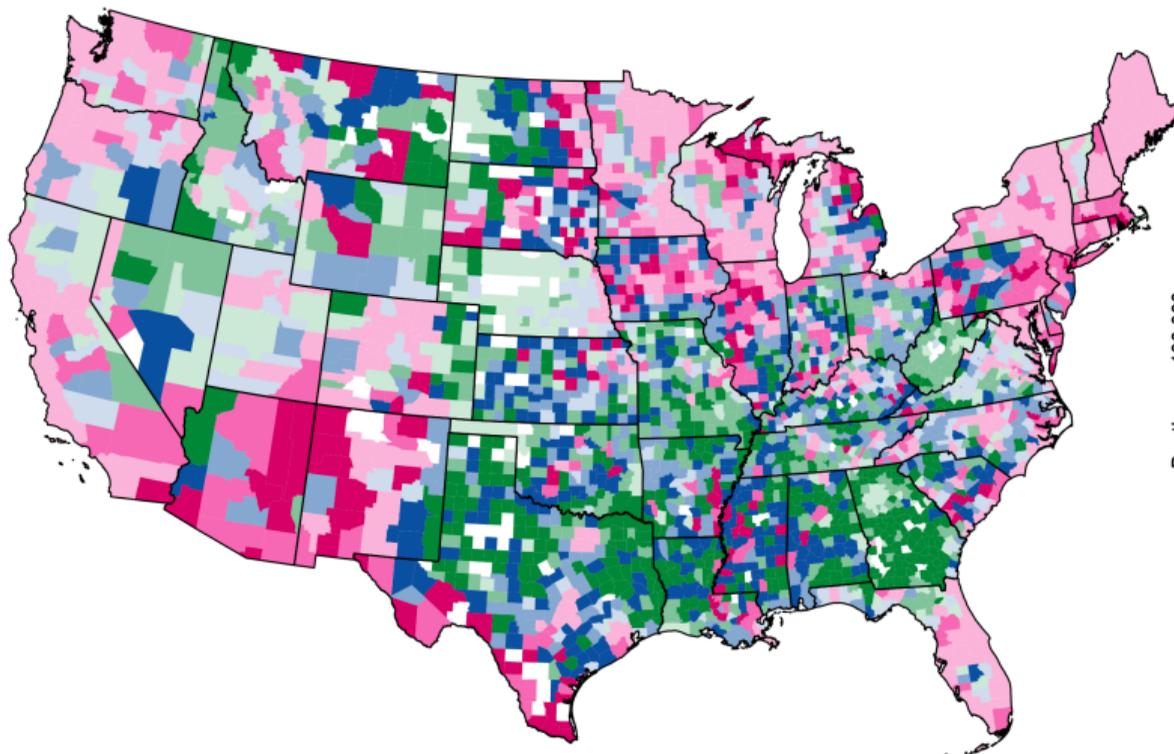
The following maps illustrate the bivariate relationship between vaccination rates and cumulative deaths in waves 4, 5 and 6. Vaccines first became available during wave 3, only being widespread over the first several months of 2021.

In the maps' legend, the rightmost categories refer to the third tercile of vaccination rates. The dark green counties are those with low vaccination rates and the third tercile of cumulative deaths. The geography of those most seriously affected by the course of the virus is quite evident.

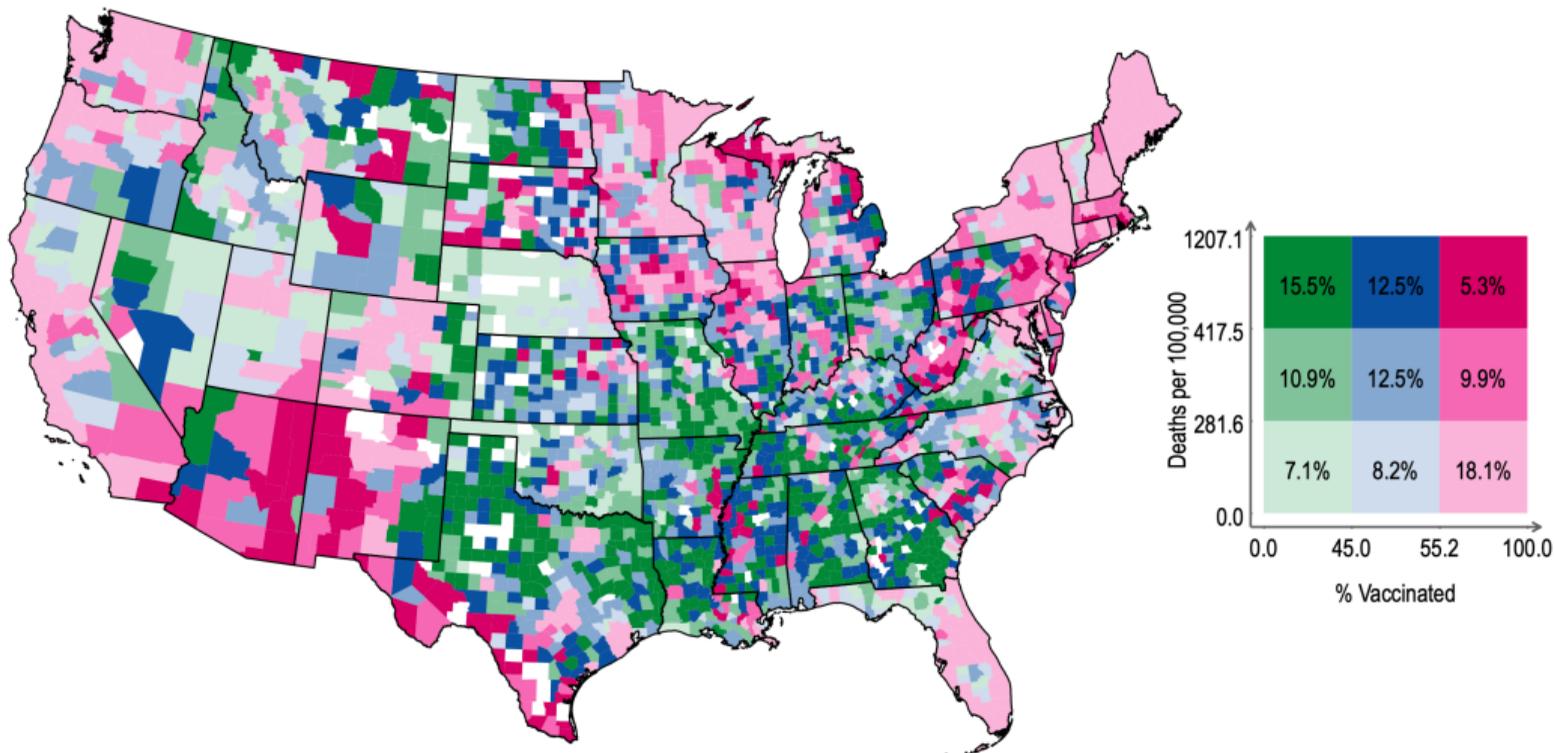
Wave 4: 1 April–31 July 2021



Wave 5: 1 August–30 November 2021



Wave 6: 1 December 2021–19 March 2022



First stage modeling

We now present the detailed econometric strategy for our investigation of these relationships among COVID-19 cases, deaths and vaccinations across time and space.

We analyze the associations between cumulative confirmed cases and deaths at the county level, each expressed per 100,000 population. Starting with wave 3, we also consider the fraction of the county population recorded as being fully vaccinated.

The models are fit separately for each of the six waves to allow for variations over those episodes in the transmissibility of the virus, the impact of vaccinations, and improvements in treatment regimes.

First stage modeling

We now present the detailed econometric strategy for our investigation of these relationships among COVID-19 cases, deaths and vaccinations across time and space.

We analyze the associations between cumulative confirmed cases and deaths at the county level, each expressed per 100,000 population. Starting with wave 3, we also consider the fraction of the county population recorded as being fully vaccinated.

The models are fit separately for each of the six waves to allow for variations over those episodes in the transmissibility of the virus, the impact of vaccinations, and improvements in treatment regimes.

In order to allow for heterogeneity within a wave, we fit mixed models (Stata's `mixed`), allowing both intercept and slopes to vary by county in this panel data context. An unstructured covariance matrix is used to provide flexibility for the random effects at the county level.

The first model for the daily county-level confirmed cases is autoregressive, with a single regressor: the county-level confirmed cases 14 days prior, capturing transmissibility of the disease. In waves 3–6, the county-level vaccination rate 14 days prior is also included.

In order to allow for heterogeneity within a wave, we fit mixed models (Stata's `mixed`), allowing both intercept and slopes to vary by county in this panel data context. An unstructured covariance matrix is used to provide flexibility for the random effects at the county level.

The first model for the daily county-level confirmed cases is autoregressive, with a single regressor: the county-level confirmed cases 14 days prior, capturing transmissibility of the disease. In waves 3–6, the county-level vaccination rate 14 days prior is also included.

The confirmed case model:

$$c_{it} = \alpha_0 + \alpha_i + \beta_0 c_{i,t-j} + \beta_i c_{i,t-j} + \gamma_0 v_{i,t-j} + \gamma_i v_{i,t-j} + \epsilon_{it}, \quad (1)$$

c_{it} and $c_{i,t-j}$ denote the cumulative confirmed cases per 100,000 in county i at time t and $t-j$, respectively. In turn, $v_{i,t-j}$ indicates the percentage of county residents that are fully vaccinated (with a second dose of a two-dose vaccine or a dose of a single-dose vaccine); α_0 , β_0 and γ_0 denote unknown fixed parameters; α_i , β_i and γ_i denote county-level random effects; and ϵ_{it} is the disturbance term.

The second model associates the daily county-level deaths with a single regressor: the county-level confirmed cases 14 days prior, capturing the mortality risk for those infected. In waves 3–6, the county-level vaccination rate 14 days prior is also included.

$$d_{it} = \delta_0 + \delta_i + \kappa_0 c_{i,t-j} + \kappa_i c_{i,t-j} + \lambda_0 v_{i,t-j} + \lambda_i v_{i,t-j} + \varepsilon_{it}. \quad (2)$$

δ_0 , κ_0 and λ_0 are unknown fixed parameters; δ_i , κ_i and λ_i denote county-level random effects; and ε_{it} is the disturbance term.

Following the estimation of the first stage models, the county-level random slopes for the lagged case regressor are predicted and added to the fixed coefficient for that variable.

Although the mean of county-level random effects is zero over the entire sample, it varies considerably at the county level for each wave, reflecting the heterogeneity in these dynamic relationships that arises from state-level and county-level characteristics and policies.

Following the estimation of the first stage models, the county-level random slopes for the lagged case regressor are predicted and added to the fixed coefficient for that variable.

Although the mean of county-level random effects is zero over the entire sample, it varies considerably at the county level for each wave, reflecting the heterogeneity in these dynamic relationships that arises from state-level and county-level characteristics and policies.

Second stage modeling

In the second stage, the outcome variables are the cross-sectional coefficients computed for each county and wave. An extensive set of fixed factors are considered as possible drivers of the transmissibility coefficients (from the case equation) and mortality risk coefficients (from the death equation).

As the dependent variable is an estimated parameter, there is a need to account for the uncertainty surrounding the estimated coefficients in the first stage models: $\hat{\beta}_i$ and $\hat{\gamma}_i$ in the case rate model (eq. 1) and $\hat{\kappa}_i$ and $\hat{\lambda}_i$ in the death rate model (eq. 2). We employ weighted least squared estimation using the precision of the first stage coefficient estimates as weights for both dependent and independent variables. The weights are applied separately for each wave's equation, rather than as a fixed set of weights for the entire estimation.

Second stage modeling

In the second stage, the outcome variables are the cross-sectional coefficients computed for each county and wave. An extensive set of fixed factors are considered as possible drivers of the transmissibility coefficients (from the case equation) and mortality risk coefficients (from the death equation).

As the dependent variable is an estimated parameter, there is a need to account for the uncertainty surrounding the estimated coefficients in the first stage models: $\hat{\beta}_i$ and $\hat{\gamma}_i$ in the case rate model (eq. 1) and $\hat{\kappa}_i$ and $\hat{\lambda}_i$ in the death rate model (eq. 2). We employ weighted least squared estimation using the precision of the first stage coefficient estimates as weights for both dependent and independent variables. The weights are applied separately for each wave's equation, rather than as a fixed set of weights for the entire estimation.

The socioeconomic and demographic factors include:

- Age and sex distribution
- Racial/ethnic distribution
- Log median household income
- High school, college completion
- Poverty rate
- Rate of owner-occupied housing

The socioeconomic and demographic factors include:

- Age and sex distribution
- Racial/ethnic distribution
- Log median household income
- High school, college completion
- Poverty rate
- Rate of owner-occupied housing

The socioeconomic and demographic factors include:

- Age and sex distribution
- Racial/ethnic distribution
- Log median household income
- High school, college completion
- Poverty rate
- Rate of owner-occupied housing

The socioeconomic and demographic factors include:

- Age and sex distribution
- Racial/ethnic distribution
- Log median household income
- High school, college completion
- Poverty rate
- Rate of owner-occupied housing

The socioeconomic and demographic factors include:

- Age and sex distribution
- Racial/ethnic distribution
- Log median household income
- High school, college completion
- Poverty rate
- Rate of owner-occupied housing

The socioeconomic and demographic factors include:

- Age and sex distribution
- Racial/ethnic distribution
- Log median household income
- High school, college completion
- Poverty rate
- Rate of owner-occupied housing

The health factors include:

- Health insurance coverage
- Diabetes prevalence
- Smoking prevalence
- Years of life expectancy
- Index of severe COVID-19 health risk (z-score)
- Access to exercise opportunities
- ICU beds per capita
- Medicaid expansion status

The health factors include:

- Health insurance coverage
- Diabetes prevalence
- Smoking prevalence
- Years of life expectancy
- Index of severe COVID-19 health risk (z-score)
- Access to exercise opportunities
- ICU beds per capita
- Medicaid expansion status

The health factors include:

- Health insurance coverage
- Diabetes prevalence
- Smoking prevalence
- Years of life expectancy
- Index of severe COVID-19 health risk (z-score)
- Access to exercise opportunities
- ICU beds per capita
- Medicaid expansion status

The health factors include:

- Health insurance coverage
- Diabetes prevalence
- Smoking prevalence
- Years of life expectancy
- Index of severe COVID-19 health risk (z-score)
- Access to exercise opportunities
- ICU beds per capita
- Medicaid expansion status

The health factors include:

- Health insurance coverage
- Diabetes prevalence
- Smoking prevalence
- Years of life expectancy
- Index of severe COVID-19 health risk (z-score)
- Access to exercise opportunities
- ICU beds per capita
- Medicaid expansion status

The health factors include:

- Health insurance coverage
- Diabetes prevalence
- Smoking prevalence
- Years of life expectancy
- Index of severe COVID-19 health risk (z-score)
- Access to exercise opportunities
- ICU beds per capita
- Medicaid expansion status

The health factors include:

- Health insurance coverage
- Diabetes prevalence
- Smoking prevalence
- Years of life expectancy
- Index of severe COVID-19 health risk (z-score)
- Access to exercise opportunities
- ICU beds per capita
- Medicaid expansion status

The health factors include:

- Health insurance coverage
- Diabetes prevalence
- Smoking prevalence
- Years of life expectancy
- Index of severe COVID-19 health risk (z-score)
- Access to exercise opportunities
- ICU beds per capita
- Medicaid expansion status

The other factors include:

- Average seasonal temperatures in summer and winter
- Average seasonal relative humidity in summer and winter
- Level of $PM_{2.5}$ pollutants
- Log population density
- 2020 presidential election Democratic vote share
- Social Vulnerability Index

The other factors include:

- Average seasonal temperatures in summer and winter
- Average seasonal relative humidity in summer and winter
- Level of $PM_{2.5}$ pollutants
- Log population density
- 2020 presidential election Democratic vote share
- Social Vulnerability Index

The other factors include:

- Average seasonal temperatures in summer and winter
- Average seasonal relative humidity in summer and winter
- Level of $PM_{2.5}$ pollutants
- Log population density
- 2020 presidential election Democratic vote share
- Social Vulnerability Index

The other factors include:

- Average seasonal temperatures in summer and winter
- Average seasonal relative humidity in summer and winter
- Level of $PM_{2.5}$ pollutants
- Log population density
- 2020 presidential election Democratic vote share
- Social Vulnerability Index

The other factors include:

- Average seasonal temperatures in summer and winter
- Average seasonal relative humidity in summer and winter
- Level of $PM_{2.5}$ pollutants
- Log population density
- 2020 presidential election Democratic vote share
- Social Vulnerability Index

The other factors include:

- Average seasonal temperatures in summer and winter
- Average seasonal relative humidity in summer and winter
- Level of $PM_{2.5}$ pollutants
- Log population density
- 2020 presidential election Democratic vote share
- Social Vulnerability Index

In the second stage empirical analysis, a regression equation is fit separately to the coefficients generated for each of the six waves. In the initial specification, all factors are included in each equation.

These estimates are refined by applying the one-covariate-at-a-time (OCMT) variable selection algorithm proposed by Chudik, Kapetanios, and Pesaran (2018) and implemented in a community-contributed routine, `ocmt`, provided by Núñez and Otero (2020), available from the SSC Archive.

In the second stage empirical analysis, a regression equation is fit separately to the coefficients generated for each of the six waves. In the initial specification, all factors are included in each equation.

These estimates are refined by applying the one-covariate-at-a-time (OCMT) variable selection algorithm proposed by Chudik, Kapetanios, and Pesaran (2018) and implemented in a community-contributed routine, `ocmt`, provided by Núñez and Otero (2020), available from the SSC Archive.

OCMT serves as an alternative approach to penalized regression for variable selection in high-dimensional linear regression models. Its objective is to find a set of predictors that is sufficient to approximate the true data generating process underlying the variable of interest. Among the several advantages of OCMT over penalized regression methods, its authors highlight ease of interpretation, its relation to classical statistical analysis, computational speed, and good performance in small samples.

As the name implies, OCMT tests the statistical significance of all covariates one at a time and selects those whose t -statistics are in absolute value greater than a given critical value. The critical value is computed using the critical value function $c_p(K, \theta) = \Phi^{-1}\left(1 - \frac{p}{2f(K, \theta)}\right)$, where $\Phi^{-1}(\cdot)$ is the inverse of the standard normal distribution function, $f(K, \theta) = cK^\theta$ for some positive constant $c = 1$ and θ , the critical value exponent; $0 < p < 1$ is the nominal size of the individual test statistics; and K is the number of covariates in the regression model of interest. All of the covariates that satisfy the stated condition are selected jointly to form the initial specification of the model.

In a second stage, OCMT uses this initial specification and once again tests the statistical significance of the covariates not selected before one at a time. The procedure continues until there are no more statistically significant covariates. OCMT is fast because the number of covariates bounds the number of stages required for convergence.

As the name implies, OCMT tests the statistical significance of all covariates one at a time and selects those whose t -statistics are in absolute value greater than a given critical value. The critical value is computed using the critical value function $c_p(K, \theta) = \Phi^{-1}\left(1 - \frac{p}{2f(K, \theta)}\right)$, where $\Phi^{-1}(\cdot)$ is the inverse of the standard normal distribution function, $f(K, \theta) = cK^\theta$ for some positive constant $c = 1$ and θ , the critical value exponent; $0 < p < 1$ is the nominal size of the individual test statistics; and K is the number of covariates in the regression model of interest. All of the covariates that satisfy the stated condition are selected jointly to form the initial specification of the model.

In a second stage, OCMT uses this initial specification and once again tests the statistical significance of the covariates not selected before one at a time. The procedure continues until there are no more statistically significant covariates. OCMT is fast because the number of covariates bounds the number of stages required for convergence.

The application of OCMT to the equation for each of the six waves provides a more parsimonious specification in which only selected factors are included in the equation. The OCMT estimates for each wave produce a different set of factors for those six episodes, with varying coefficients for those factors selected for multiple waves.

These cross-section estimates can be considered as a system of six equations with differing specifications and time-varying coefficients. This system of equations is then estimated with a novel application of Zellner's seemingly unrelated regression (SUR) estimator: Stata's `sureg`. The usual context for SUR is a set of OLS equations for several units (firms, industries, countries) in a time-series context. In our application of SUR, the equations correspond to different time periods, and the observations are the 3,014 U.S. counties in our analysis.

The application of OCMT to the equation for each of the six waves provides a more parsimonious specification in which only selected factors are included in the equation. The OCMT estimates for each wave produce a different set of factors for those six episodes, with varying coefficients for those factors selected for multiple waves.

These cross-section estimates can be considered as a system of six equations with differing specifications and time-varying coefficients. This system of equations is then estimated with a novel application of Zellner's seemingly unrelated regression (SUR) estimator: Stata's `sureg`. The usual context for SUR is a set of OLS equations for several units (firms, industries, countries) in a time-series context. In our application of SUR, the equations correspond to different time periods, and the observations are the 3,014 U.S. counties in our analysis.

The usual rationale for SUR as a systems estimator is the degree to which each equation's error process might be contemporaneously correlated with other units' errors at each *point in time*. If those correlations are sizable, SUR can yield efficiency gains relative to single-equation estimation of each equation.

In our context, the error correlations that can be exploited are those *for each county* over the six waves of the pandemic. Those correlations should be sizable, as they reflect unobservable factors at the county level that have not been captured by the time-invariant regressors selected for each wave. The degree to which these correlations increase the precision of the estimates is evaluated by the Breusch–Pagan test for independence, computed with the `sureg` option `corr`, with the null hypothesis that the 6x6 residual correlation matrix is diagonal. Under the null, this test statistic is distributed $\chi^2(m)$, where $m = 15$, the number of subdiagonal elements in the matrix. The null is strongly rejected for both applications of the SUR technique.

The usual rationale for SUR as a systems estimator is the degree to which each equation's error process might be contemporaneously correlated with other units' errors at each *point in time*. If those correlations are sizable, SUR can yield efficiency gains relative to single-equation estimation of each equation.

In our context, the error correlations that can be exploited are those *for each county* over the six waves of the pandemic. Those correlations should be sizable, as they reflect unobservable factors at the county level that have not been captured by the time-invariant regressors selected for each wave. The degree to which these correlations increase the precision of the estimates is evaluated by the Breusch–Pagan test for independence, computed with the sureg option `corr`, with the null hypothesis that the 6x6 residual correlation matrix is diagonal. Under the null, this test statistic is distributed $\chi^2(m)$, where $m = 15$, the number of subdiagonal elements in the matrix. The null is strongly rejected for both applications of the SUR technique.

Table: OCMT SUR cross-section results for $\hat{\beta}_i$ in cumulative case model (eq. 1) with $j = 14$

Wave starting:	3/20	7/20	10/20	4/21	8/21	12/21
Age 1-19 yrs (%)	0.00643**			0.01688***		
Age 20-39 yrs (%)	-0.00407		-0.00610***	0.00915*		
Age 40-59 yrs (%)			-0.00092		-0.00271	-0.00382***
Age 60-79 yrs (%)	0.00201		-0.00228*	0.00307		
Black (%)	-0.00266***	-0.00105*		0.00036	-0.00287***	-0.00085***
Hispanic (%)	-0.00021	-0.00177**		-0.00693***	-0.00036	-0.00069***
Male (%)	0.00007		-0.00695***	-0.00838*	-0.00105	-0.00464***
Median income (log)		-0.01190	0.04332		0.09835***	0.05423***
Social vulnerability index	0.06349**	0.04213	0.04354**	0.29297***	0.11183***	
HS completion (%)	-0.19051	0.16884	0.21472***	0.80663***	0.26546***	
Some college (%)		-0.04406	-0.09222**		-0.15706***	
Poverty rate (%)	0.28348**	-0.08389	-0.08792	-0.74218***		-0.11775**
Owner-occupied housing (%)	-0.17082**		-0.02006	0.12142		
Medicaid expansion	-0.02986***	-0.04012**	-0.00963	0.22587***	-0.04535***	-0.02412***
Uninsured (%)	-0.00115	-0.00398*	-0.00119	-0.00040	-0.00172*	-0.00556***
Diabetes rate (%)		0.00268	0.00632***	-0.00545*	0.00125	0.00116*
Smoking (%)	-0.66763***		0.12001	2.17326***	0.44826***	0.92696***
Life expectancy (years)		-0.00170	-0.00039	0.02006***	-0.00094	0.00256***
Health risk index	0.00048	-0.02014*	-0.03158***	-0.02133	0.01646***	0.00230
Access to exercise (%)	-0.00032	0.13441***	-0.00072		0.04302**	0.01037
ICU (beds per 100K)	-0.00008		-0.00001	0.00055	-0.00002	0.00011
Observations	3041					

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

Table: OCMT SUR cross-section results for $\hat{\beta}_i$ in cumulative case model (eq. 1) with $j = 14$

Wave starting:	3/20	7/20	10/20	4/21	8/21	12/21
IPop. density (log)	0.02034***		0.02436***	0.04162***	0.02732***	0.00946***
$PM_{2.5}$		-0.00888**	0.00611***	-0.02018***	0.00025	0.00175
Summer avg. temp. (C)	-0.00072	0.01394***	-0.00719***	0.01869***	-0.00383**	
Summer rel. hum. (%)	-0.00114*			-0.00318**		-0.00076***
Winter avg. temp. (C)	0.00473***	-0.01075***	-0.00101	0.03565***	-0.00215**	0.00553***
Winter rel. hum. (%)	-0.00103	0.00186		0.00582***	0.00112	0.00180***
Democratic share 2020 (%)	0.11693**		0.08423***		0.10194***	0.09820***
Constant	0.73270***	0.04286	0.49819	-3.78922***	-0.81472***	-0.34931*
Observations	3041					

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

For the cumulative confirmed case model the OCMT selection process identifies several regressors that are relevant for each wave: for instance, an indicator of the state's Medicaid expansion and the percent of residents uninsured. Both of those factors have negative effects on the transmissibility of the virus in each wave, but the magnitude of those effects varies considerably across waves.

A number of other factors are identified as relevant in only certain waves. The log of population density has a significant positive effect in all but wave 2, the last pre-vaccination wave. The Breusch–Pagan test for the relevance of correlations among residuals from each wave rejects its null hypothesis with a p-value of 0.00.

For the cumulative confirmed case model the OCMT selection process identifies several regressors that are relevant for each wave: for instance, an indicator of the state's Medicaid expansion and the percent of residents uninsured. Both of those factors have negative effects on the transmissibility of the virus in each wave, but the magnitude of those effects varies considerably across waves.

A number of other factors are identified as relevant in only certain waves. The log of population density has a significant positive effect in all but wave 2, the last pre-vaccination wave. The Breusch–Pagan test for the relevance of correlations among residuals from each wave rejects its null hypothesis with a p-value of 0.00.

Table: OCMT SUR cross-section results for $\hat{\kappa}_i$ in cumulative death model (eq. 2) with $j = 14$

Wave starting:	3/20	7/20	10/20	4/21	8/21	12/21
Age 1-19 yrs (%)			-0.00100***			
Age 20-39 yrs (%)		-0.00022**	-0.00121***			-0.00003
Age 40-59 yrs (%)	0.00035*		-0.00123***			
Age 60-79 yrs (%)		0.00010	-0.00109***			
Black (%)	0.00005*				0.00002	0.00001**
Male (%)	-0.00044***	-0.00021*				-0.00005*
Median income (log)		-0.00029	-0.00382***	0.00342	-0.00172	
Social vulnerability index				0.01357**	0.00179*	
HS completion (%)				0.06545***	0.00926**	
Some college (%)				-0.00792	-0.00344	
Poverty rate (%)				-0.00873	0.00058	
Owner-occupied housing (%)						0.00316***
Uninsured (%)	-0.00021**			0.00059**	0.00015***	
Diabetes rate (%)		0.00011		0.00031	0.00009	0.00003**
Smoking (%)				0.02940	-0.02284***	
Life expectancy (years)		-0.00034***		-0.00005	-0.00018**	
Health risk index		-0.00040		0.00123	0.00062**	
Access to exercise (%)	-0.00274			-0.00308		
ICU (beds per 100K)						0.00001***
Constant	0.05068***	0.05094***	0.16734***	-0.12993	0.04313***	0.00155
Observations	3041					

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

Table: OCMT SUR cross-section results for $\hat{\kappa}_j$ in cumulative death model (eq. 2) with $j = 14$

Wave starting:	3/20	7/20	10/20	4/21	8/21	12/21
Pop. density (log)	0.00149***		0.00062***			
$PM_{2.5}$	-0.00024				0.00030***	0.00008**
Summer avg. temp. (C)			-0.00030***	0.00115***	-0.00034***	
Summer rel. hum. (%)						0.00000
Winter avg. temp. (C)		0.00012***	-0.00001	-0.00036	0.00003	
Winter rel. hum. (%)	-0.00027***					
Democratic share 2020 (%)	0.00592*					
Constant	0.05068***	0.05094***	0.16734***	-0.12993	0.04313***	0.00155
Observations	3041					

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

For the cumulative deaths model, the OCMT selection process considers far fewer factors as influencing the mortality risk from infection in waves 1–3, corresponding to the period from March 2020–March 2021. In the following two waves (April–November 2021) several additional factors are identified as having important effects, while few factors appear in wave 6 (December 2021–March 2022), perhaps reflecting the improved treatments now available.

As a robustness test, we recomputed the second stage estimates using LASSO to implement variable selection. The results were qualitatively similar, with OCMT generally retaining fewer regressors than the LASSO technique.

Summary

This study analyzed two years of daily data on COVID-19 cases and deaths at the U.S. county level. The two-stage modeling approach allows for unobservable factors to affect both the estimated transmissibility of the virus and the mortality risk for those infected, treating each of six distinct waves of the pandemic.

The cross-sectional coefficients produced in the first stage can then be used to identify the sociodemographic, health, climate, pollution and political factors that have played important roles in these outcomes, allowing for variations in model specification and coefficients over the six waves. This flexible approach provides considerable insight to the process by which the course of the pandemic has been affected over time and space.

Summary

This study analyzed two years of daily data on COVID-19 cases and deaths at the U.S. county level. The two-stage modeling approach allows for unobservable factors to affect both the estimated transmissibility of the virus and the mortality risk for those infected, treating each of six distinct waves of the pandemic.

The cross-sectional coefficients produced in the first stage can then be used to identify the sociodemographic, health, climate, pollution and political factors that have played important roles in these outcomes, allowing for variations in model specification and coefficients over the six waves. This flexible approach provides considerable insight to the process by which the course of the pandemic has been affected over time and space.

References I

-  Allcott, H. et al. (2020). *What Explains Temporal and Geographic Variation in the Early US Coronavirus Pandemic?* NBER Working Paper No. 27965.
https://www.nber.org/system/files/working_papers/w27965/w27965.pdf.
-  Aron, Janine and John Muellbauer (June 2022). “Excess Mortality Versus COVID-19 Death Rates: A Spatial Analysis of Socioeconomic Disparities and Political Allegiance Across U.S. States”. In: *Review of Income and Wealth* 68.2, pp. 348–392. doi: 10.1111/roiw.12570.
url: <https://ideas.repec.org/a/bla/revinw/v68y2022i2p348-392.html>.
-  Baum, Christopher F. and Miguel Henry (2022). “Socio-economic and demographic factors influencing the spatial spread of COVID-19 in the USA”. In: *International Journal of Computational Economics and Econometrics* 12.4, pp. 366–380. url: <https://ideas.repec.org/a/ids/ijcome/v12y2022i4p366-380.html>.

References II

-  Brown, Caitlin S. and Martin Ravallion (July 2020). *Inequality and the Coronavirus: Socioeconomic Covariates of Behavioral Responses and Viral Outcomes Across US Counties*. NBER Working Papers 27549. National Bureau of Economic Research, Inc. url: <https://ideas.repec.org/p/nbr/nberwo/27549.html>.
-  Carozzi, Felipe, Sandro Provenzano, and Sefi Roth (2022). “Urban Density and COVID-19: understanding the US experience”. In: *Annals of Regional Science* 28, pp. 1–32. doi: [10.1007/s00168-022-01193-z](https://doi.org/10.1007/s00168-022-01193-z).
-  Chudik, Alexander, George Kapetanios, and M Hashem Pesaran (2018). “A one covariate at a time, multiple testing approach to variable selection in high-dimensional linear regression models”. In: *Econometrica* 86.4, pp. 1479–1512.
-  Desmet, K. and R. Wacziarg (2022). “JUE Insight: Understanding spatial variation in COVID-19 across the United States”. In: *Journal of Urban Economics* 127. doi: [10.1016/j.jue.2021.103332](https://doi.org/10.1016/j.jue.2021.103332).

References III

-  Doti, James L. (2021). “Examining the impact of socioeconomic variables on COVID-19 death rates at the state level”. In: *Journal of Bioeconomics* 23, pp. 15–53. doi: 10.1007/s10818-021-09309-9.
-  Jones, Bradley (2022). *The Changing Political Geography of COVID-19 Over the Last Two Years*. Tech. rep. Version 2022-03-03. Available at <https://www.pewresearch.org/politics/2022/03/03/the-changing-political-geography-of-covid-19-over-the-last-two-years/>. Pew Research Center.
-  Knittel, Christopher R. and Bora Ozaltun (June 2020). *What Does and Does Not Correlate with COVID-19 Death Rates*. NBER Working Papers 27391. National Bureau of Economic Research, Inc. url: <https://ideas.repec.org/p/nbr/nberwo/27391.html>.

References IV

-  McLaren, John (July 2021). “Racial Disparity in COVID-19 Deaths: Seeking Economic Roots with Census Data”. In: *The B.E. Journal of Economic Analysis & Policy* 21.3, pp. 897–919. doi: 10.1515/bejeap-2020-0371. url: <https://ideas.repec.org/a/bpj/bejeap/v21y2021i3p897-919n12.html>.
-  Mukherji, Nivedita (2022). “The social and economic factors underlying the incidence of COVID-19 cases and deaths in US counties during the initial outbreak phase”. In: *Review of Regional Studies* 52, pp. 127–150.
-  Naqvi, Asjad (2022). *BIMAP: Stata module to produce bivariate maps*. Statistical Software Components, Boston College Department of Economics. url: <https://ideas.repec.org/c/boc/bocode/s459063.html>.

References V

-  Núñez, Héctor M. and Jesús Otero (2020). *OCMT: Stata module to perform multiple testing approach in high-dimensional linear regression*. Statistical Software Components, Boston College Department of Economics. url: <https://ideas.repec.org/c/boc/bocode/s458850.html>.
-  Papageorge, Nicholas W. et al. (June 2020). *Socio-Demographic Factors Associated with Self-Protecting Behavior during the Covid-19 Pandemic*. NBER Working Papers 27378. National Bureau of Economic Research, Inc. url: <https://ideas.repec.org/p/nbr/nberwo/27378.html>.
-  Welsch, David (2022). “The Impact of Mask Usage on COVID-19 Deaths: Evidence from US Counties Using a Quasi-Experimental Approach”. In: *The B.E. Journal of Economic Analysis & Policy* 22.1, pp. 1–28. doi: 10.1515/bejeap-2021-0157.

Appendix: Spatial modeling

A number of studies have relied on spatial modeling, taking into account the spread of a communicable disease across adjacent geographic units. These include the analysis by Allcott et al. (2020), who consider policy interventions and social distancing in the context of a model of disease transmissions, and find that population and density were the key factors in the first six months of the pandemic.

Carozzi, Provenzano, and Roth (2022) study the impact of urban population density as it affected the timing of outbreaks. Desmet and Wacziarg (2022) also consider the early stages of the pandemic by evaluating county-level correlates of cases and deaths. Baum and Henry (2022) use a spatial autoregressive model to analyze the spread of the pandemic over the first 14 months, taking county-level demographic factors into account.

Appendix: Spatial modeling

A number of studies have relied on spatial modeling, taking into account the spread of a communicable disease across adjacent geographic units. These include the analysis by Allcott et al. (2020), who consider policy interventions and social distancing in the context of a model of disease transmissions, and find that population and density were the key factors in the first six months of the pandemic.

Carozzi, Provenzano, and Roth (2022) study the impact of urban population density as it affected the timing of outbreaks. Desmet and Wacziarg (2022) also consider the early stages of the pandemic by evaluating county-level correlates of cases and deaths. Baum and Henry (2022) use a spatial autoregressive model to analyze the spread of the pandemic over the first 14 months, taking county-level demographic factors into account.

Appendix: Demographic and socioeconomic studies

Many studies have considered the importance of demographic and socioeconomic factors on the evolution of cases and deaths. These include the analysis by Papageorge et al. (2020), evaluating survey data from the first months of the pandemic and individuals' behavioral changes such as social distancing and mask wearing.

Mukherji (2022) evaluates the importance of socioeconomic factors in the initial outbreak of the pandemic, developing a social vulnerability index. Brown and Ravallion (2020) focus on socioeconomic measures of inequality and how they interact with the spread of the disease in the first half of 2020. McLaren (2021) studies the roots of racial disparities in pandemic deaths.

Appendix: Demographic and socioeconomic studies

Many studies have considered the importance of demographic and socioeconomic factors on the evolution of cases and deaths. These include the analysis by Papageorge et al. (2020), evaluating survey data from the first months of the pandemic and individuals' behavioral changes such as social distancing and mask wearing.

Mukherji (2022) evaluates the importance of socioeconomic factors in the initial outbreak of the pandemic, developing a social vulnerability index. Brown and Ravallion (2020) focus on socioeconomic measures of inequality and how they interact with the spread of the disease in the first half of 2020. McLaren (2021) studies the roots of racial disparities in pandemic deaths.

Knittel and Ozaltun (2020) find that some common factors are not correlated with death rates in the early phase of the pandemic. Doti (2021) finds that statewide mandates became more effective in preventing deaths in late 2020. Welsch (2022) evaluates linkages between mask wearing and pandemic deaths, taking the potential endogeneity into account. Aron and Muellbauer (2022) study excess ('all-causes') mortality through February 2021, linking socioeconomic disparities and political factors to state-level deaths. They find that political factors played an important role in reducing mortality disparities.

Most of the cited literature addresses the early phases of the pandemic. In our study, we analyze the pandemic's evolution over a longer period.

Knittel and Ozaltun (2020) find that some common factors are not correlated with death rates in the early phase of the pandemic. Doti (2021) finds that statewide mandates became more effective in preventing deaths in late 2020. Welsch (2022) evaluates linkages between mask wearing and pandemic deaths, taking the potential endogeneity into account. Aron and Muellbauer (2022) study excess ('all-causes') mortality through February 2021, linking socioeconomic disparities and political factors to state-level deaths. They find that political factors played an important role in reducing mortality disparities.

Most of the cited literature addresses the early phases of the pandemic. In our study, we analyze the pandemic's evolution over a longer period.

Appendix: Relationships among regressors

To illustrate some of the relationships among these variables, we present bivariate maps of several pairs of factors that play an important role in the second stage analysis:

- log median income and population density
- poverty and lack of health insurance
- life expectancy and health risk index
- Democratic vote share 2020 and Medicaid expansion

Appendix: Relationships among regressors

To illustrate some of the relationships among these variables, we present bivariate maps of several pairs of factors that play an important role in the second stage analysis:

- log median income and population density
- poverty and lack of health insurance
- life expectancy and health risk index
- Democratic vote share 2020 and Medicaid expansion

Appendix: Relationships among regressors

To illustrate some of the relationships among these variables, we present bivariate maps of several pairs of factors that play an important role in the second stage analysis:

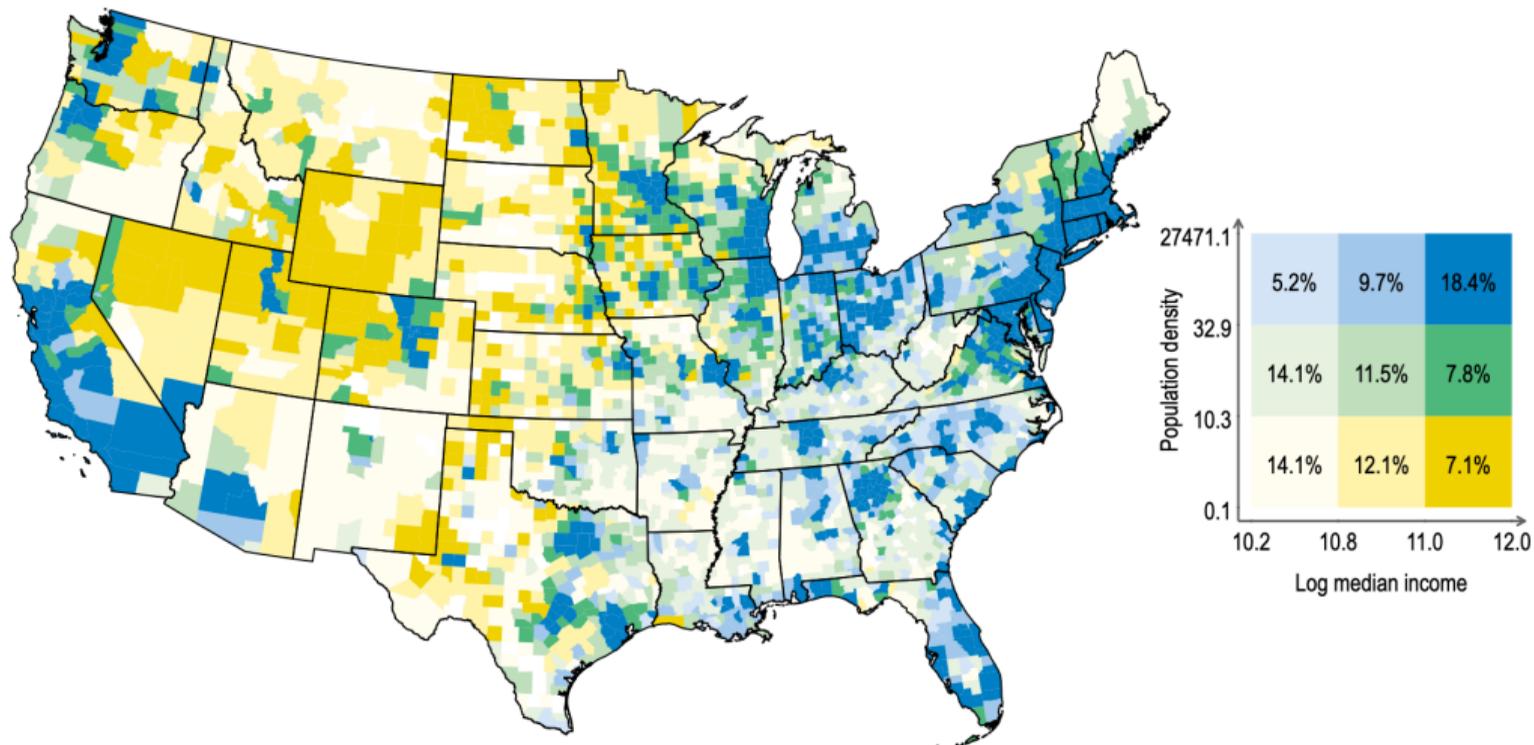
- log median income and population density
- poverty and lack of health insurance
- life expectancy and health risk index
- Democratic vote share 2020 and Medicaid expansion

Appendix: Relationships among regressors

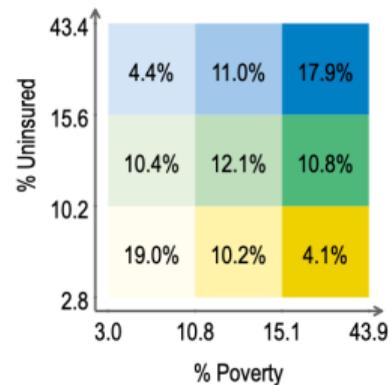
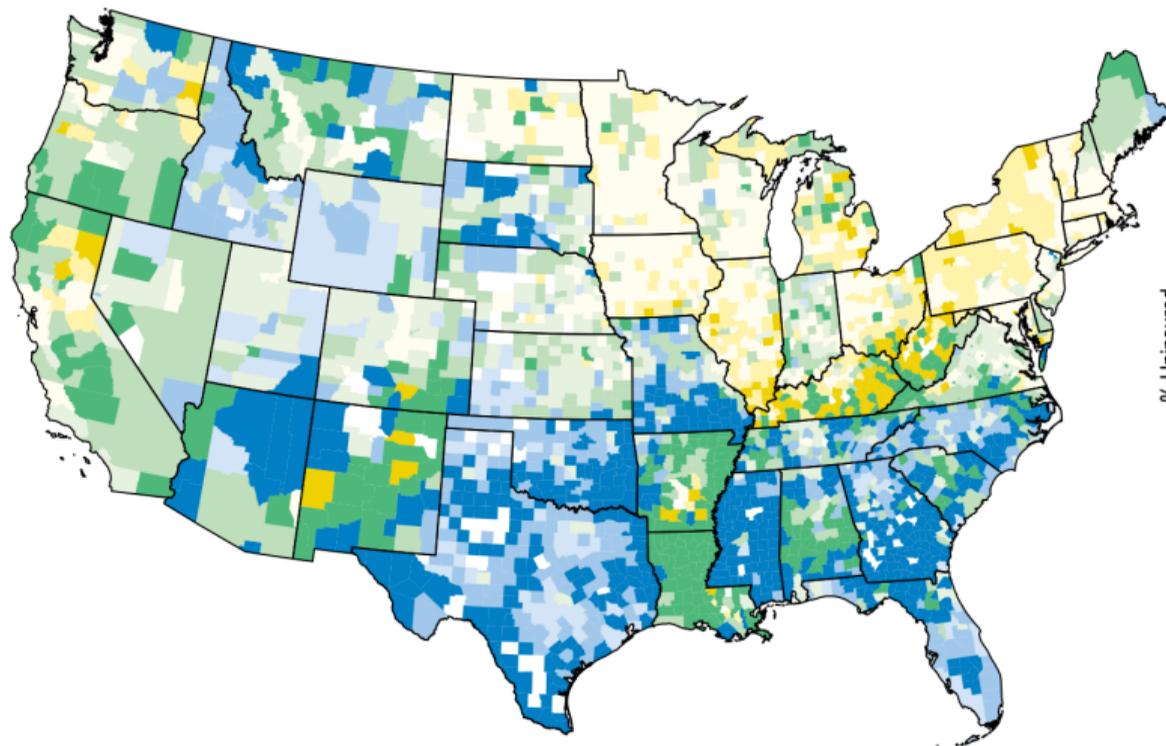
To illustrate some of the relationships among these variables, we present bivariate maps of several pairs of factors that play an important role in the second stage analysis:

- log median income and population density
- poverty and lack of health insurance
- life expectancy and health risk index
- Democratic vote share 2020 and Medicaid expansion

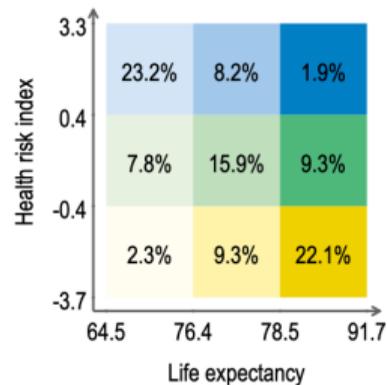
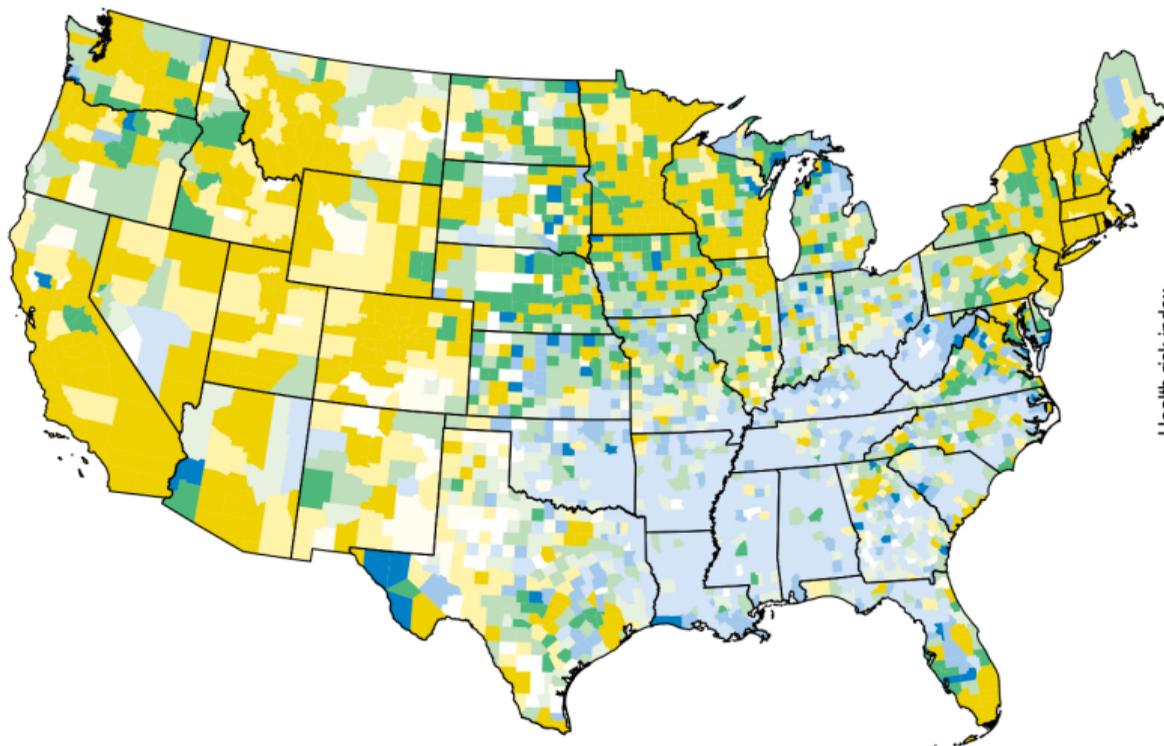
log median income and population density



poverty and lack of health insurance



life expectancy and health risk index



Democratic vote share 2020 and Medicaid expansion

