

Mask Mandates Can Rapidly and Efficiently Limit COVID-19 Spread: Month-Over-Month Effectiveness of Governmental Policies in Reducing the Number of New COVID-19 Cases in 37 US States and the District of Columbia

Michael J. Maloney,
Proof School, San Francisco, CA

Nathaniel J. Rhodes, Pharm.D., M.Sc.,
Chicago College of Pharmacy, and the
Pharmacometrics Center of Excellence, Midwestern University

and

Paul R. Yarnold, Ph.D.
Optimal Data Analysis, LLC

SARS-CoV-2 is the beta-coronavirus responsible for COVID-19. Facemask use has been qualitatively associated with reduced COVID-19 cases, but no study has quantitatively assessed the impact of government mask mandates (*MM*) on new COVID-19 cases across multiple US States. We used a non-parametric machine-learning algorithm to test the *a priori* hypothesis that *MM* were associated with reductions in new COVID-19 cases. Publicly available data were used to analyze new COVID-19 cases from 37 States and the District of Columbia (i.e., “38 States”). We conducted confirmatory All-States and State-Wise analyses, validity analyses [e.g., leave-one-out (LOO) and bootstrap resampling], and covariate analyses. No statistically significant difference in the daily number of new COVID-19 infections was discernable in the All-States analysis. In State-Wise LOO validity analysis, 11 States exhibited reductions in new COVID-19, and reductions in four of these States (AK, MA, MN, VA) were statistically significant in bootstrap

resampling. Only the Social Capital Index predicted *MM* success (training $p < 0.028$ and LOO $p < 0.013$). Results obtained when studying the impact of *MM* on COVID-19 cases varies as a function of the heterogeneity of the sample being considered, providing clear evidence of *Simpson's Paradox* and thus of confounded findings. As such, studies of *MM* effectiveness should be conducted on disaggregated data. Since transmissions occur at the individual rather than at the collective level, additional work is needed to identify optimal social, psychological, environmental, and educational factors which will reduce the spread of SARS-CoV-2 and facilitate *MM* effectiveness across diverse settings.

The significance of respiratory droplets and airborne viral particles in the spread of SARS-CoV-2 and propagation of the COVID-19 pandemic gained early attention.¹ Asymptomatic carriage of SARS-CoV-2 is important to community spread because these individuals can transmit the virus in exhaled breath.^{2,3} Corollary research reported that face masks are protective⁴⁻⁷ and may limit the severity of COVID-19 among individuals who become infected.⁸ Indeed, the Centers for Disease Control and Prevention (CDC) has advocated wearing a facemask when social isolation is impossible or impractical,⁹ and recently reported that dining-out is among the riskiest known activities during the coronavirus pandemic—since face masks *are not used* when people are eating and drinking.¹⁰ Other activities associated with an increased risk of SARS-CoV-2 transmission include singing¹¹ and aerosol generating procedures in the healthcare setting.¹² Small droplet formation, which can lead to aerosolization of the virus, is a major driver of infection risk,¹³ thus wearing a facemask is expected to reduce inhalation of aerosolized virus.⁴ Indeed, low rates of positive serology were observed (i.e., $< 4\%$, $n=1/27$) when the healthcare workers followed aerosol minimizing procedures *and* used personal protective equipment to perform deep respiratory sampling in patients with COVID-19.¹⁴ Thus, there is consensus in the

scientific, clinical, and business communities that appropriate wearing of facemasks is a “best-practice” personal behavior, which can reduce the chances of being infected by, or of infecting others with, SARS-CoV-2 and other airborne disease-causing microbes.

Mathematical and time-series models of COVID-19 data collected in the US investigated salutary effects of masking on reducing spread of the pandemic.^{15,16} These models suggest that a consistent application of best-practice public-health orders and citizen behaviors is needed to slow the spread of the pandemic.¹⁷ A hybrid modeling approach used by the Institute for Health Metrics and Evaluation predicted that if 95% of people in the US wore masks outside their homes, the number of projected deaths from COVID-19 would drop by half within four months.¹⁸ Thus, experts hypothesize that Mask Mandate (*MM*) orders should result in reduced COVID-19 transmission and fatalities—assuming that the individuals in the community adhere to the *MM*.

Encouraged by research findings which showed that wearing masks reduces the number of new COVID-19 cases, and motivated to intervene by increasing numbers of cases, the Governors of 37 US States and the District of Columbia (hereafter called the “38 States”) issued *MM* orders. However, the lack of clear, consistent Global, Federal, and State guidance

on the use of face coverings (e.g., mandates *vs.* recommendations), and on who should be using masks (e.g., essential workers *vs.* private citizens), has initiated a nation-wide naturalistic experiment.¹⁵ In the absence of clear guidance and of direct leadership, variations in *MM* acceptance and enforcement have arisen within and between US States.^{17, 19-22}

Promulgating and exacerbating confusion regarding the effectiveness of wearing face masks, adjudicating the effectiveness of *MM* orders has largely been left to *qualitative interpretations* and *visual/numerical inspections* of the epidemic curves depicting the number of new cases over time. Qualitative assessment cannot objectively test the *a priori* hypothesis that a *MM* reduces the incidence of COVID-19 cases. Accordingly, using publicly-available data, we employ non-parametric maximum-accuracy machine-learning to *quantitatively* test the confirmatory (*a priori*) hypothesis that imposing a *MM* reduced daily number of new COVID-19 cases in the month *after vs. before* the *MM*.

Methods

Hypotheses

First, merging data from all 38 States in which the Governor issued a *MM*, we test the *a priori* hypothesis that “*All-States*” imposition of a *MM* reduced the daily number of new COVID-19 cases in the month *after vs. in the month before* the *MM*. States were included if State leadership (e.g., Governors) had mandated the use of face-masks by all public-facing employees.

Second, considering each State individually, we test the *a priori* hypothesis that “*State-Wise*” imposition of *MM* reduced the daily number of new COVID-19 cases in the month *after vs. in the month before* the *MM*.

Data

The daily number of new COVID-19 cases was obtained separately for each State in the month *before* the *MM* was made (a maximum of 30 days), and in the month *after* the *MM* was made (a maximum of 31 days). Case reports occurring before the *MM* were dummy-coded as class=0, and case reports occurring after the *MM* were coded as class=1. Data from *The New York Times*, based on reports from state and local health agencies, were initially downloaded from GitHub on August 23, 2020 (data re-confirmed on September 23, 2020)²³ and cross referenced against the COVID-19 State Policy Database (updated May 28, 2020) for dates that *MM* were issued in each State.²⁴

Our decision to compare daily number of new cases in the month before *vs.* the month after the *MM* was based on methodological and statistical considerations. First, *MM* orders were made because the day-over-day increases in the number of new cases threatened the ability of available health care resources to meet patient and health care worker needs. Therefore, the month following the *MM* naturally offers great opportunity for increased mask wearing to meaningfully reduce the number of new cases. Second, some amount of time is required after the *MM* is made for citizens (who wish to comply) to obtain masks, learn how and when to wear them, and how to care for them, and for new infections occurring in the two-week period prior to the *MM* to develop symptoms after the *MM* was made. Third, to maximize the available statistical power of our analyses for estimating short-term effects (see *Statistical Analysis*), we extracted the number of cases for up to 30 days before and up to 31 days after the *MM* was implemented (yielding a maximum of n=61 data points for each state).

Additional attributes, obtained for individual States, were used in an attempt to discriminate States for which a *MM* did vs. did not reduce the number of new COVID-19 cases. Numerical attributes included population estimates,²⁵ gross domestic product (State and per capita),²⁶ age (with arbitrary brackets of 0-18, 18-25, 26-34, 35-54, 55-64, and 65+),²⁷ pre-pandemic (January 2020) number of unemployed persons and rate,²⁸ homelessness,²⁹ shelter beds,³⁰ incarceration number and rate,³¹ number of nursing home residents,³² population density (people per square mile in 2015),³³ urban overcrowding (number of houses having >1 person per room),³⁴ severe urban overcrowding (number of houses having >1.5 people per room),³⁴ ethnicity (white, black, and Hispanic categories),³⁵ Social Capital Index (a measure of social cohesion, where higher positive numbers indicate greater societal cohesion),³⁶ and obesity.³⁷ Categorical attributes included if the State's Governor wore a mask in public after the *MM* was made (0=Did not wear mask, 1=wore mask) which was assessed using media reports and public database image searches, the Governor's political party (0=Republican, 1=Democrat),³⁸ and the Governor's gender (0=Female, 1=Male).³⁹

Statistical Analysis

We used confirmatory Optimal Discriminant (or Data) Analysis (ODA) to ascertain the *predictive*

accuracy attained by using maximum-accuracy models to evaluate the *hypothesized reduction in the number of new cases* occurring subsequent to the *MM*.⁴⁰⁻⁴² ODA is the moniker of the statistically-motivated non-parametric machine-learning algorithm which identifies the cut-point for an ordered *attribute* (i.e., "independent variable"), or the assignment rule for a categorical attribute, which yields the *strongest achievable accuracy* in discriminating between two or more class categories (i.e., "dependent variable") for a given sample and hypothesis.⁴⁰⁻⁴² In the present study all hypotheses were confirmatory.⁴¹ For both the All-States (38-State) and the State-Wise analysis, the *a priori* hypothesis was that the *MM* (up to 31 days) reduced the number of new cases.^{43,44}

An ODA model is obtained by iterating through every possible different assignment rule consistent with the *a priori* hypothesis, and identifying the model achieving the highest effect strength for sensitivity (ESS) statistic.^{40,41} An index of classification accuracy adjusted to remove the effect of chance, ESS is a function of the mean sensitivity of the model achieved across class categories, standardized such that 0 represents the discriminatory accuracy which is expected by chance, 100 represents errorless (perfect) discrimination, and values of $-100 \leq \text{ESS} < 0$ represent accuracy which is worse than expected by chance. For a two-category *class variable* such as studied herein (i.e., month before *MM*=class 0 vs. month after *MM*=class 1),

$$\text{Mean Classification Accuracy (\%)} = 100 \times [(\text{sensitivity for class 0} + \text{sensitivity for class 1})/2] \quad (1)$$

and

$$\text{ESS} = 100 \times [(\text{Mean Classification Accuracy} - 50) / 50]. \quad (2)$$

A nonparametric permutation test which requires no distributional assumptions is used to assess the statistical significance (p value) of the achieved ESS.^{40, 43, 44} Point estimates of p values for tests of statistical hypotheses reported herein are indicated as being statistically significant ($p \leq 0.05$) at either the *experimentwise* (confirmatory $p_{observed} \leq p_{Sidak\ adjusted\ Bonferroni}$) or the *per-comparison* (confirmatory $p \leq 0.05$) criterion.⁴¹

Internal model validation was conducted using bootstrap resampling^{45,46} ($n=25,000$) of the State-Wise models within the ODA package for R.⁴⁷ Specifically, resampling with 50% replacement was conducted using the LOO (i.e., validity) ODA model confusion matrix for each State. Exact discrete 95% confidence intervals (CI) for the given ODA model (“Model”) and for randomly scrambled observations from the model (“Chance”) were obtained.⁴⁵ Graphical depictions of new cases organized across days, and histograms of bootstrapped 95% CIs, were created using *ggplot2* for R.⁴⁸

Potential generalizability of the ODA model applied to classify an independent random sample is assessed using a one-sample “leave-one-out” (LOO) jackknife analysis.⁴¹ The sample size offers greater than 90% power to identify an ODA model of moderate strength (i.e., $25 \leq ESS < 50$).⁴⁹

Finally, we used exploratory ODA analyses to model the association of inter-State differences in a host of variables (e.g., age, wealth, population density) conceivably capable of discriminating States which did vs. did not exhibit the predicted decline in number of new cases in the month following a MM.⁴¹

Results

All-States Analysis

The *mean* number of new COVID-19 cases before vs. after the MM for the sample of 38 States was 654 ($N=1138$, $SD=1357$) vs. 639 ($N=1177$,

$SD=975$), respectively. In three instances, the daily report of new cases appeared to be outliers (GA 4/12/20, LA 5/29/20, and VA 5/6/20) and were coded as missing (i.e., coded as -999). After exclusion of outliers, case numbers in the confirmatory evaluation matched the original data download for all states except Oregon ($n=12/61$) and Virginia ($n=5/12$), which were within $\pm 6\%$ of the original data: these numerical changes in data values (none of which were near the optimal model cut-point) did not change the ODA model, or the accuracy (ESS) or statistical significance (p) of the model in training or LOO analysis. Confirmatory ODA was used to test the *a priori* hypothesis that imposition of a MM reduced the *number* of new COVID-19 cases (the attribute) one-month *after vs.* one-month *before* the MM (the class variable). The resulting confirmatory ODA model was: if >4 new cases-per-day then predict pre-MM; otherwise if ≤ 4 new cases-per-day, then predict post-MM. In training (using all observations from the sample) analysis, this model correctly classified 98.1% of the pre-MM cases, and 6.1% of the post-MM cases: $ESS=4.1$ (a weak effect), $p < 0.13$.

Thus, when considering all 38 States having a MM as a single sample, there was no statistically significant difference in the daily number of new COVID-19 infections in the month before vs. the month after the MM.

State-Wise Analysis

In Table 1 the *mean* number of new daily cases is numerically *lower* post Mask Mandate (MM) in 14 of 38 States (indicated in red font).

Table 1: State-Wise Descriptive Statistics for Number of New Cases Pre- vs. Post-MM

State	Before MM		After MM	
	Mean	SD	Mean	SD
AK	9.8	4.8	2.4	2.0
AL	229.4	64.1	390.3	158.0
AR	93.7	75.4	205.0	114.8

AZ	245.6	88.2	547.7	355.1
CA	1417.9	390.6	2147.9	464.8
CO	338.6	81.2	421.1	148.5
DC	68.0	52.3	155.3	52.7
DE	147.2	100.1	153.7	69.2
FL	754.3	263.8	863.7	305.2
GA	704.1	315.6	653.4	218.1
HI	17.1	9.4	3.4	4.3
IL	1564.2	504.7	2182.9	668.9
IN	530.7	134.4	559.0	132.2
KY	159.5	87.9	178.8	85.7
LA	760.2	606.1	399.2	237.0
MA	1925.7	718.4	1041.5	660.5
MD	385.2	292.1	911.0	240.9
ME	26.4	12.6	39.7	16.7
MI	1147.5	344.6	572.2	255.6
MN	637.7	116.4	383.6	101.5
MS	222.9	68.5	277.7	81.8
NE	179.6	152.3	296.1	109.6
NH	59.3	17.4	80.8	29.6
NJ	1480.3	1528.0	2936.7	808.8
NM	117.1	47.0	146.3	51.3
NV	120.8	38.8	124.9	42.7
NY	7479.5	2836.8	4170.2	2206.4
OH	503.5	288.7	574.1	106.9
OR	61.0	14.9	59.9	29.0
PA	1051.8	605.8	1150.4	355.1
RI	138.1	120.7	278.0	87.1
TX	923.8	222.6	1266.6	380.9
UT	65.9	56.0	137.8	34.7
VA	933.2	243.9	656.0	230.7
VT	25.1	16.2	5.5	5.2
WA	299.2	132.8	252.7	108.5
WV	31.8	17.7	28.4	19.5
WY	14.6	12.9	11.1	6.6

Table 2 gives findings of ODA used to test the *a priori* hypothesis that State-Wise imposition of a *MM* (made by 38 different State Governors) *reduced the number of new COVID-19 cases in the month after vs. in the month before the MM*—with each State considered

individually. In each analysis the class variable was a binary dummy-coded variable signifying the period before *vs.* after the *MM*, and the attribute was the daily number of new cases.

In Table 2 States are ordered from the strongest to weakest ESS of the confirmatory ODA model in training analysis. Among all 38 States, training data for Minnesota (MN, listed first in Table 2) offer greatest support for the *a priori* hypothesis. The ODA model is: if >480 new cases-per-day, predict pre-*MM*; otherwise if ≤480 cases-per-day, predict post-*MM*. This model correctly classified 28 of 30 (93.3%) of the pre-*MM* days, and 27 of 31 (87.1%) of the post-*MM* days, thereby yielding a *strong*⁴¹ ESS=80.4, *p*<0.0001. However, this model showed modest instability in LOO jackknife analysis (*relatively strong*⁴¹ ESS=73.9, *p*<0.0001). In contrast, data for Alaska (AK, second row in Table 2) had the second-strongest training effect (*strong* ESS=80.3, *p*<0.0001), which was stable in LOO analysis. Among all 38 States, only AK returned a *strong* confirmation of the *a priori* hypothesis in training and LOO analysis.

Of the total of 38 States in this study, 11 (23.7%) achieved confirmatory *p*<0.05 at the experimentwise criterion, and two (5.3%) achieved confirmatory *p*<0.05 at the per-comparison criterion⁴¹ (Table 2). Three States (MN, AK, HI) yielded a training model indicating a *strong* confirmation⁴¹ (i.e., ESS≥75) of the *a priori* hypothesis. Five States (MA, VT, VA, MI, NY) yielded a training model indicating a *relatively strong* confirmation⁴¹ (50≤ESS<75). And, three States (LA, WA, WY) yielded a training model indicating a *moderate* confirmation⁴¹ (25≤ESS<50) of the *a priori* hypothesis.

Except for GA and WY, States having a numerically greater *mean* number of new daily cases pre- *vs.* post-*MM* (Table 1) also had statistically-significant (at either the experimentwise or the per-comparison criterion⁴²) confirmatory (Table 2) and LOO (Table 3) ODA models supporting the *a priori* hypothesis.

Table 2: Maximum-Accuracy Confirmatory Models of State-Wise Mask Mandate (MM) Effectiveness

State	Pre-MM	Sensitivity (%)	Post-MM	Specificity (%)	Training $p <$	ESS	LOO $p <$	ESS
MN	>480 Cases	28/30 (93.3)	≤480 Cases	27/31 (87.1)	0.0001	80.4	0.0001	73.9
AK	>4 Cases	27/30 (90.0)	≤4 Cases	28/31 (90.3)	0.0001	80.3	0.0001	80.3
HI	>6 Cases	26/30 (86.7)	≤6 Cases	28/31 (90.3)	0.0001	77.0	0.0001	64.0
MA	>1174 Cases	29/30 (96.7)	≤1174 Cases	23/31 (74.2)	0.0001	70.9	0.0001	64.3
VT	>14 Cases	23/30 (76.7)	≤14 Cases	29/31 (93.6)	0.0001	70.2	0.0001	66.8
VA	>691 Cases	27/29 (93.1)	≤691 Cases	23/31 (74.2)	0.0001	67.3	0.0001	67.3
MI	>732 Cases	26/30 (86.7)	≤732 Cases	24/31 (77.4)	0.0001	64.1	0.0001	50.8
NY	>4762 Cases	26/30 (86.7)	≤4762 Cases	22/31 (71.0)	0.0001	57.6	0.0001	51.0
LA	>340 Cases	26/30 (86.7)	≤340 Cases	17/30 (56.7)	0.0027	43.3	0.0037	36.7
WA	>294 Cases	17/30 (56.7)	≤294 Cases	25/31 (80.6)	0.011	37.3	0.0028	37.3
WY	>15 Cases	13/30 (43.3)	≤15 Cases	27/31 (87.1)	0.036	30.4	0.04	24.0
PA	>1389 Cases	12/30 (40.0)	≤1389 Cases	26/31 (83.4)	0.16	23.9	0.07	20.5
OR	>41 Cases	28/30 (93.3)	≤41 Cases	8/31 (25.8)	0.25	19.1	0.23	9.5
GA	>886 Cases	8/29 (27.6)	≤886 Cases	28/31 (90.3)	0.34	17.9	0.21	11.0
WV	>36 Cases	11/30 (36.7)	≤36 Cases	25/31 (80.6)	0.31	17.3	0.18	14.0
DE	>203 Cases	6/30 (20.0)	≤203 Cases	28/31 (90.3)	0.65	10.3	0.34	7.0
OH	>768 Cases	4/30 (13.3)	≤768 Cases	30/31 (96.8)	0.68	10.1	0.30	6.8
KY	>64 Cases	28/30 (93.3)	≤64 Cases	5/31 (16.1)	0.7	9.5	0.30	9.5
IN	>396 Cases	27/30 (90.0)	≤396 Cases	6/31 (19.4)	0.72	9.4	0.38	6.1
NV	>189 Cases	3/30 (10.0)	≤189 Cases	30/31 (96.8)	0.81	6.8	0.99	-77.1
NE	>455 Cases	3/30 (10.0)	≤455 Cases	30/31 (96.8)	0.81	6.8	0.82	-3.0
AL	>107 Cases	30/30 (100)	≤107 Cases	2/31 (6.4)	0.84	6.4	0.52	3.1
NJ	>4296 Cases	1/30 (3.33)	≤4296 Cases	31/31 (100)	0.94	3.3	0.88	-3.1
TX	>477 Cases	30/30 (100.0)	≤477 Cases	1/31 (3.2)	0.96	3.2	0.51	3.2
FL	>534 Cases	25/30 (83.3)	≤534 Cases	6/31 (19.4)	0.97	2.7	0.77	-3.9
NM	>216 Cases	2/30 (6.7)	≤216 Cases	29/31 (93.6)	0.98	0.2	0.94	-6.3
MS	>138 Cases	28/30 (93.3)	≤138 Cases	2/31 (6.4)	0.99	-0.2	0.95	-6.8
UT	>176 Cases	2/30 (6.7)	≤176 Cases	28/31 (90.3)	0.99	-3.0	0.99	-12.8
DC	>224 Cases	1/30 (3.3)	≤224 Cases	29/31 (93.6)	0.99	-3.1	0.88	-3.1
MD	>1154 Cases	1/30 (3.3)	≤1154 Cases	29/31 (93.6)	0.99	-3.1	0.94	-6.3
CO	>210 Cases	28/30 (93.3)	≤210 Cases	1/31 (3.2)	0.99	-3.4	0.89	-3.4
NH	>34 Cases	28/30 (93.3)	≤34 Cases	1/31 (3.23)	0.99	-3.4	0.95	-6.8
ME	>61 Cases	1/30 (3.3)	≤61 Cases	27/31 (87.1)	0.99	-9.6	0.99	-93.4
IL	>1133 Cases	26/30 (86.7)	≤1133 Cases	1/31 (3.2)	0.99	-10.1	0.98	-10.1
AR	>34 Cases	24/30 (80.0)	≤34 Cases	1/31 (3.2)	0.99	-16.8	0.99	-20.1
AZ	>161 Cases	24/30 (80.0)	≤161 Cases	1/31 (3.2)	0.99	-16.8	0.99	-20.1
RI	>273 Cases	9/30 (30.0)	≤273 Cases	16/31 (51.6)	0.99	-18.4	0.99	-41.7
CA	>2315 Cases	2/30 (6.7)	≤2315 Cases	21/31 (67.7)	0.99	-25.6	0.99	-32.2

For all 11 States in Table 2 with training $p < 0.05$, Table 3 presents 95% exact discrete CIs obtained for model and for chance.⁴⁵ CIs were estimated by bootstrap resampling of the ODA model confusion matrix from each respective State obtained in LOO analysis. For a given State, if the upper bound (UB) of the CI for chance falls *beneath* the lower bound (LB) of the CI for the model, this column indicates NO 95% CI overlap—so the confirmatory hypothesis is *accepted* (i.e., the *MM* reduced the number of new cases). And, for a given State, if the UB of the CI for chance falls *above* the LB of the CI for the model, this column indicates YES 95% CI overlap—so the confirmatory hypothesis is

rejected (i.e., the *MM* did not reduce the number of new cases). In other words, if the exact discrete CIs for model and chance *do not overlap*, then the ODA model identified in LOO analysis is considered statistically robust and the findings are likely to be generalizable.

As seen, for AK, VA, MA, and MN, the *lower bound* of the 95% exact discrete CI *for model* is greater than the *upper bound* of the 95% exact discrete CI *for chance*—thus the result is statistically significant. For the remaining seven States, the upper bound for chance exceeded the lower bound for the model, so therefore the result was not statistically significant in bootstrap validity analysis.

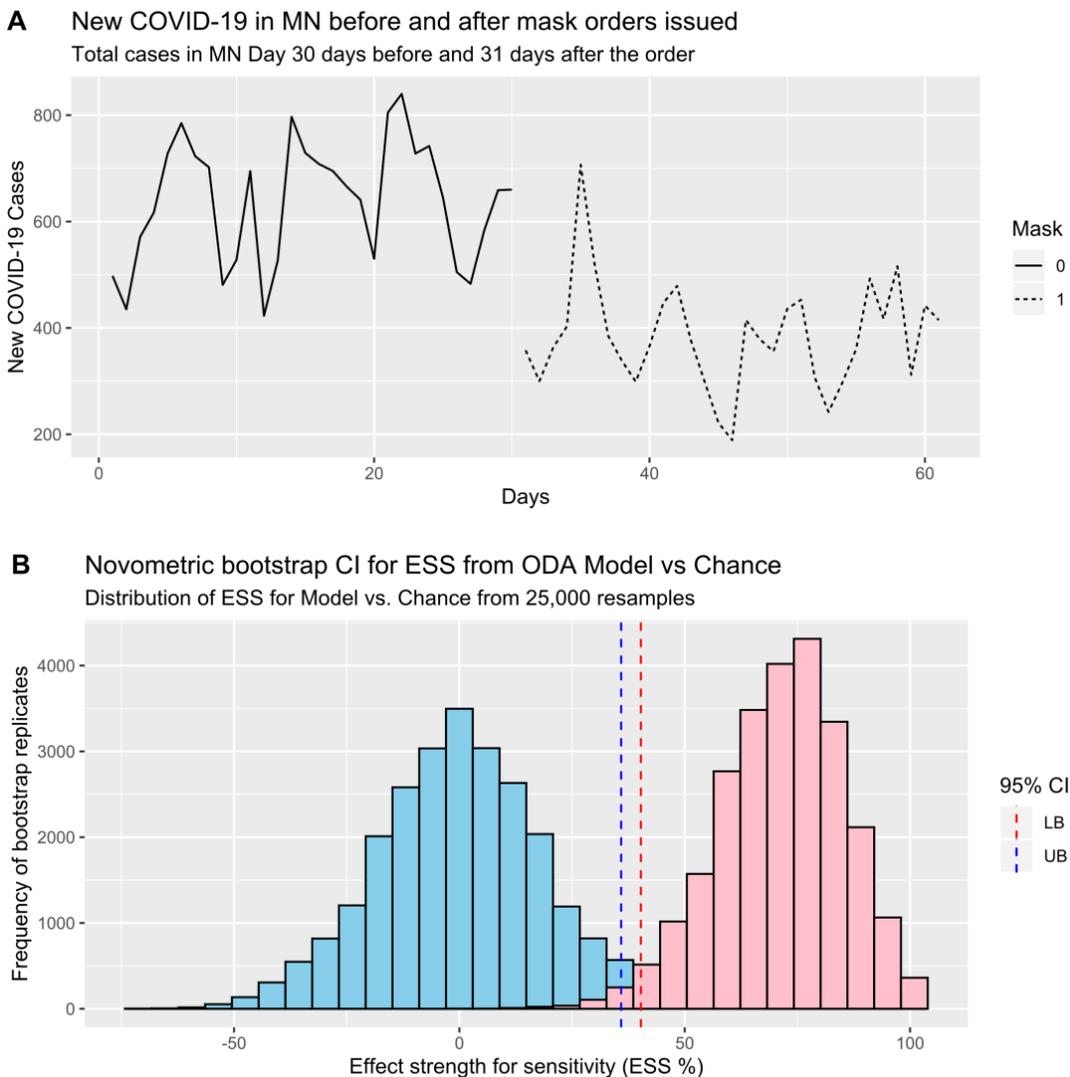
Table 3: LOO Confirmatory p , and Bootstrap 95% CI for LOO ESS, for 11 States Demonstrating Mask Mandate (*MM*) Effectiveness vis-à-vis Training $p < 0.05$

State	$p <$	Chance 95% CI	Model 95% CI	95% CI overlap
AK	0.0001	-36.1 – 35.9	64.6 – 100	No
VA	0.0001	-35.8 – 35.9	49.9 – 100	No
MA	0.0001	-35.0 – 34.8	41.7 – 92.9	No
MN	0.0001	-35.9 – 36.1	40.3 – 95.2	No
HI	0.0001	-36.1 – 35.9	34.9 – 91.7	Yes
VT	0.0001	-35.7 – 34.9	30.3 – 86.7	Yes
NY	0.0001	-36.1 – 35.8	25.5 – 86.1	Yes
LA	0.0037	-34.8 – 35.7	25.0 – 84.6	Yes
MI	0.0001	-35.8 – 36.1	11.1 – 76.6	Yes
WA	0.0028	-34.4 – 35.0	0.00 – 68.4	Yes
WY	0.04	-33.3 – 33.5	-10.1 – 58.3	Yes

The daily number of new COVID-19 cases one-month pre- and one-month post-*MM* is seen in Figure 1.A for MN. Figure 1.B is a graphical summary of model and chance bootstrap results ($n=25,000$ resamples with 50% replacement). In both Figures the blue histogram represents *chance*, the red histogram represents *the model*, the Y axis is frequency of the result in the bootstrap analysis, and the X axis is train-

ing ESS. Dashed vertical lines inside the Figure represent the lower 95% CI for the model (red dashes), and the upper 95% CI for chance (blue dashes). When the blue dashes are on the same side as the blue histogram, and red dashes are on the same side as the red histogram, the effect is statistically significant. Evaluating MN in training analysis, the *a priori* hypothesis is accepted with experimentwise $p < 0.05$.

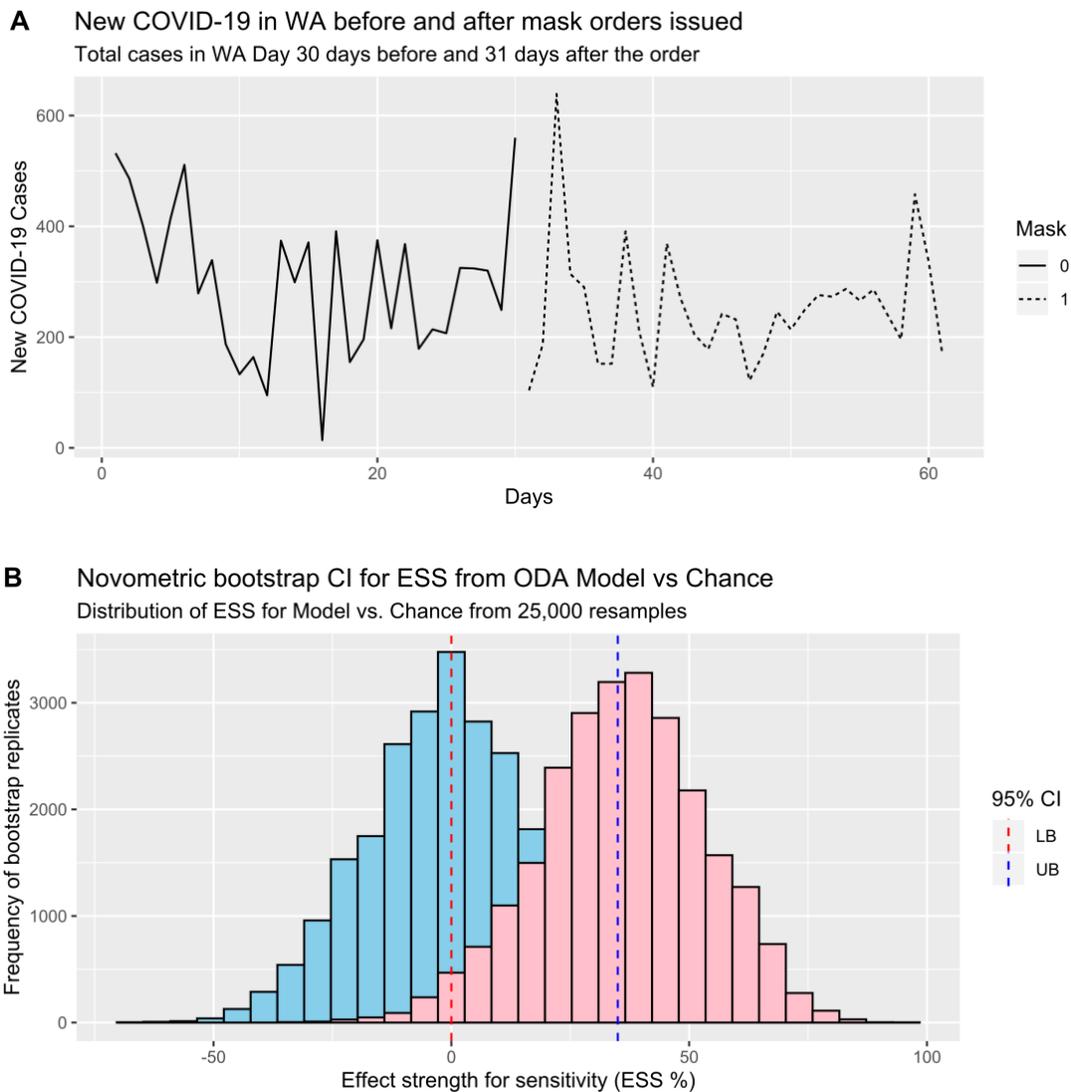
Figure 1. Representative state (MN) with a strong impact of Mask Mandate (*MM*) order



In contrast, Figure 2.A presents the daily number of new COVID-19 cases one-month pre- and one-month post-*MM* for WA, and Figure 2.B is a graphical summary of the model and chance bootstrap results. Because the *blue*

dashes are on the side of the *red* histogram, and the red dashes are on the side of the *blue* histogram, the hypothesized effect is *not* statistically significant for WA.

Figure 2. Representative state (WA) with a moderate impact of Mask Mandate (*MM*) order



Modeling Success vs. Failure of *MM* Order in Reducing Number of New COVID-19 Cases

Finally, ODA was used to discriminate the four States (AK, MA, MN, VA) having a significantly lower number of new cases after the *MM* vs. the 34 States which did not support the *a priori* hypothesis. The only attribute which met the per-comparison criterion for statistical

significance⁴¹ was SCI, a measure of societal cohesion and the closeness of social ties.⁵⁰ The ODA model was: if $SCI \leq 0.311$, then predict the *MM* did not reduce the number of new cases of COVID-19; otherwise if $SCI > 0.311$ then predict the *MM* reduced the number of new

cases. In training ($p < 0.028$) and LOO ($p < 0.013$) analysis this model correctly classified all four States in which the *MM* reduced the number of new cases, and 25 of 35 (71.43%) States in which the *MM* failed to reduce the number of new cases. This level of classification accuracy corresponds to $ESS = 71.43$, indicating a relatively strong effect.⁴¹

Discussion

An *All-States* analysis was conducted using merged data for 38 US States. Results indicated *no statistically significant decrease* in the number of new daily COVID-19 cases in the month before *vs.* after the *MM* was made. In contradistinction, a *State-Wise* analysis identified *clear evidence of State-specific MM effectiveness in reducing the number of new cases within one month of enactment* for multiple States (Tables 2 and 3). Such clear inconsistency is attributable to improper analysis of merged samples creating an analytic anomaly that is known as *Simpson's Paradox*—which may identify spurious relationships, may fail to identify actual relationships, and/or may over- or under-weight the strength of effects—in every area of quantitative empirical science.^{41,46} *Paradoxical confounding* exists when the finding obtained by statistical analysis conducted for a combined sample is different than the finding obtained by the same analysis conducted separately for the constituent samples.⁵¹ In the setting of a public health emergency, such as the current pandemic, accurate knowledge of the effectiveness of public health interventions is essential so that public health authorities can communicate this to the public.

Combining disparate groups is the primary operating procedure in government, scientific, and media discussions about the number of new and cumulative daily, monthly, and total number of COVID-19 cases and fatalities. Examples of such combined groups include global (the maximum-possible combined group), regional (e.g., the northern hemisphere),

continental (e.g., Europe, Africa), and country (e.g., China, US, India, Brazil, France, UK, Israel) daily and cumulative case counts.⁵² However, disparities in regional increases in numbers of new cases prompted some analysts to issue a warning against focusing on National trends.⁵³ Findings presented herein—that four of 38 (10.5%) States reliably demonstrated the hypothesized reduction in number of new cases in the month after the *MM* was made—similarly suggest a warning against focusing on State trends. Cognizance of the need to combine groups in a manner that prevents invalid statistical conclusions begs the question of how small must a “catchment area” be to count new cases so as to prohibit paradoxical confounding. Unfortunately, it has been shown that Simpson's Paradox can even occur for single-subject, “N-of-1” designs.⁵⁴ Thus, while accurate case counts and definitions are important for tracking the epidemiology of the pandemic, such measures do not necessarily reflect the “reality on the ground” in any specific community, nor can they reliably discern the individual effectiveness of public health interventions.

The current methodology of conducting after-the-fact pandemic accounting for convenient and familiar (statistically unmotivated) groups is obviously inadequate when it is considered from a problem-solving perspective that focuses on getting (and staying) ahead of the pandemic curve. Others have argued, early in the pandemic, that rapid and comprehensive contact tracing will empower understanding of the situations and the processes underlying new COVID-19 infections.⁵⁵ As is the case in fighting a wildfire, progress is made by going to the active location each day, extinguishing blazes, and reducing available fuel, rather than by assessing the number of acres which burned the day before. However, contact tracing is most effective when community spread is low, so multifaceted interventions are needed to control the spread of the pandemic and presently wear-

ing a mask is the most efficacious known means of slowing the pandemic.^{5, 6, 8, 15, 16}

We found that SCI, a measure of social closeness and ties,⁵⁰ was a relatively strong predictor of *MM* effectiveness. This finding is notable in that SCI could be a facilitating factor that makes mandated public health interventions more efficient. For example, populations with higher levels of social cohesion are inherently more likely to work together when faced with a common problem based on a shared sense of identity and social bonds.⁵⁶ Variation in social cohesion may partially explain why *MM* was so very clearly effective in some States and not as clearly effective in others.⁵⁷ This finding suggests a corollary hypothesis that States having a lower SCI may require different and multifaceted interventions to optimally address the pandemic. Supportive evidence for this was recently observed in a CDC analysis of *MM* effectiveness in different counties in Kansas.⁵⁸

Limitations of This Study

Our analytical approach has several limitations. Temporal counts of the number of new cases are an imprecise measure—including the month-over-month methodology we used here—for multiple reasons. Even with perfect *MM* compliance, there is an up to 14-day lag time between infection and presentation of symptoms, and a 20-day lag between infection and presentation at a healthcare facility. Thus, after a *MM* order is made, a time-lag may exist to achieve maximal effect. Our approach was focused on testing the immediate impact of *MM*, and future time-series analyses are warranted. *MM* compliance is imperfect, and enforcement is difficult. Indeed, individual *MM* compliance is expected to be heterogeneous based on experimental evidence from behavioral psychology. Individuals are susceptible to cognitive biases and often use heuristics to address complex problems when uncertainty is high. Examples of such biases and heuristics related to *MM* use

include availability bias (e.g., knowing someone who is sick vs. not), anchoring bias (e.g., fixation on initial recommendations to avoid mask use to preserve PPE), and substitution heuristics (e.g., physical discomfort with wearing a mask is substituted for mask efficacy).⁵⁹ The impact of these and other biases and heuristics on mask effectiveness is currently understudied.

Because transmissions occur at the level of the individual, rather than at a State or National level, governments should emphasize best-practice personal behaviors, adequate contact tracing, and minimization of high-risk exposures. Indeed, we posit that the efficacy of mask wearing is not at question as this depends only on the use of the mask. Instead, the effectiveness of *MM* can be eroded by multiple factors including asymptomatic carriage, willful ignorance of *MM* and lack of knowledge of how and when to wear masks. Moreover, our State-Wise analysis supports the *a priori* hypothesis that *MM do indeed have a significant salutary effect on the spread of the pandemic*, but more work is needed to define the social, psychological, environmental, and educational facilitators that maximize *MM* effectiveness.

Experimentally speaking, the best test of public health intervention effectiveness (e.g., *MM*)—using a method known as an “A (before intervention), B (after intervention)” design— involves simultaneous uniform application, enforcement, and education around the intervention, coupled with exacting real-time contact-tracing. A science-governed, consistent approach to disease management could reduce public misunderstanding in this realm, and diminish mistrust of other public health interventions.⁵⁸ Our analysis is limited because we were not able to utilize a rigorous AB analytic framework to test our *a priori* hypotheses. Rather, we took advantage of the naturalistic experiment created by the heterogeneous response to the pandemic across States to highlight the paradoxical results obtained when

applying an All-States followed by a State-Wise analytic approach. The present analysis focused on absolute case numbers which should not be confused with case rates; however, given the relatively short time-horizon for this study (one month before and after *MM*) it is doubtful that changes in test availability would have affected our findings. Additionally, our data sources have been updated at various time points over time, however a sensitivity analyses did not find any differences in case counts that would meaningfully change our conclusions.

Implications of This Study

This study has important *methodological* and *public health policy* implications for future research and practice.

Methodological implications of this study highlight the uncertain validity of research conclusions which are based on data analysis conducted for a combined sample consisting of observations from heterogeneous populations.

- We identified paradoxical confounding: the effect of a mask mandate (*MM*) on the number of cases differed if analyzing pooled data from 38 States *vs.* analyzing data for each State considered alone.
- Considering only National or State-level trends may thus obscure evidence of the effectiveness of important public health interventions (e.g., mask wearing), since transmission events are stochastic and occur at the level of individuals. Tracking state and county case counts and those of neighboring states with high border crossing rates appears very prudent.

- Understanding the impact of personal behavior on *MM* effectiveness, and thorough contact tracing for positive cases, are important specific factors underlying person-to-person transmission.
- Societal factors reflected by measures such as the SCI measure of cohesion may increase the efficiency of *MM* interventions in specific regions.
- Research is needed to identify interventions to reduce SARS-CoV-2 transmission *for socially diverse populations*.

Public health policy implications of this study highlight the short-term effectiveness of *MMs* in reducing the number of new COVID-19 cases, and logistical advantages of *MM* initiatives.

- This study provides quantitative demonstration of *MM* efficacy in *rapidly* reducing the number of new cases—within one month of initiating the *MM*.
- This study provides quantitative demonstration of *MM* efficacy in *substantially* reducing the number of new COVID-19 cases: AK, VA, MA, and MN all had statistically significant, relatively strong or strong reductions in the number of new COVID-19 infections within one month of initiating the *MM*.
- Jackknife validity analysis suggests that the *speed* and *efficacy* of *MMs* in these states will be observed for independent random samples of people.

References

- ¹Zhang R, Li Y, Zhang AL, Wang Y and Molina MJ (2020). Identifying airborne transmission as the dominant route for the spread of COVID-19. *Proceedings of the National Academy of Sciences*, 117, 14857-14863. doi: 10.1073/pnas.2009637117
- ²Bae SH, Shin H, Koo H-Y, Lee SW, Yang JM, et al. (2020). Asymptomatic transmission of SARS-CoV-2 on evacuation flight. *Emerging Infectious Diseases*, 26. doi: 10.3201/eid2611.203353
- ³Centers for Disease Control and Prevention (CDC) (2020). How COVID-19 spreads. See: <https://www.cdc.gov/coronavirus/2019-ncov/prevent-getting-sick/how-covid-spreads.html> Accessed on: September 21, 2020
- ⁴Leung NHL, Chu DKW, Shiu EYC, Chan K-H, McDevitt JJ, et al. (2020). Respiratory virus shedding in exhaled breath and efficacy of face masks. *Nature Medicine*, 26, 676-680. doi: 10.1038/s41591-020-0843-2
- ⁵Wang X, Ferro EG, Zhou G, Hashimoto D and Bhatt DL (2020). Association between universal masking in a health care system and SARS-CoV-2 positivity among health care workers. *JAMA*, 324, 703. doi: 10.1001/jama.2020.12897
- ⁶Clase CM, Fu EL, Joseph M, Beale RCL, Dolovich MB, et al. (2020). Cloth masks may prevent transmission of COVID-19: An evidence-based, risk-based approach. *Annals of Internal Medicine*, 173, 489-491. doi: 10.7326/m20-2567
- ⁷Mitze T, Kosfeld R, Rode J and Wälde K (2020). Face masks considerably reduce COVID-19 cases in Germany [pre-print]. medRxiv. doi: 10.1101/2020.06.21.20128181
- ⁸Gandhi M, Beyrer C and Goosby E (2020). Masks do more than protect others during COVID-19: Reducing the inoculum of SARS-CoV-2 to protect the wearer. *Journal of General Internal Medicine*, 35, 3063–3066. doi: 10.1007/s11606-020-06067-8
- ⁹Centers for Disease Control and Prevention (CDC) (2020). Considerations for wearing masks: Help slow the spread of COVID-19. See: <https://www.cdc.gov/coronavirus/2019-ncov/prevent-getting-sick/cloth-face-cover-guidance.html> Accessed on: August 7, 2020
- ¹⁰Fisher KA, Tenforde MW, Feldstein LR, Lindsell CJ, Shapiro NI, et al. (2020). Community and close contact exposures associated with COVID-19 among symptomatic adults ≥ 18 years in 11 outpatient health care facilities — United States, July 2020. *MMWR Morbidity and Mortality Weekly Report*, 69, 1258-1264. doi: 10.15585/mmwr.mm6936a5
- ¹¹Hamner L, Dubbel P, Capron I, Ross A, Jordan A, et al. (2020). High SARS-CoV-2 attack rate following exposure at a choir practice — Skagit County, Washington, March 2020. *MMWR Morbidity and Mortality Weekly Report*, 69, 606-610. doi: 10.15585/mmwr.mm6919e6
- ¹²Semple MG, Tran K, Cimon K, Severn M, Pessoa-Silva CL, et al. (2012). Aerosol generating procedures and risk of transmission of acute respiratory infections to healthcare workers: A systematic review. *PLoS ONE*, 7: e35797. doi: 10.1371/journal.pone.0035797
- ¹³Anfinrud P, Stadnytskyi V, Bax CE and Bax A (2020). Visualizing speech-generated oral fluid droplets with laser light scattering. *New England Journal of Medicine*, 382, 2061-2063. doi: 10.1056/NEJMc2007800

¹⁴Gao CA, Bailey JI, Walter JM, Coleman JM, Malsin ES, et al. (2020). Bronchoscopy on intubated COVID-19 patients is associated with low infectious risk to operators at a high-volume center using an aerosol-minimizing protocol [pre-print]. medRxiv. doi: 10.1101/2020.08.30.20177543

¹⁵Lyu W and Wehby GL (2020). Community use of face masks and COVID-19: Evidence from a natural experiment of state mandates in the US. *Health Affairs*, 39, 1419-1425. doi: 10.1377/hlthaff.2020.00818

¹⁶Eikenberry SE, Mancuso M, Iboi E, Phan T, Eikenberry K, et al. (2020). To mask or not to mask: Modeling the potential for face mask use by the general public to curtail the COVID-19 pandemic. *Infectious Disease Modelling*, 5, 293-308. doi: 10.1016/j.idm.2020.04.001

¹⁷Althouse BM, Wallace B, Case B, Scarpino SV, Berdhal A, et al. (2020). The unintended consequences of inconsistent pandemic control policies [pre-print]. medRxiv. doi: 10.1101/2020.08.21.20179473

¹⁸Institute for Health Metrics and Evaluation (2020). New IHME COVID-19 forecasts see nearly 300,000 deaths by December 1, however, consistent mask-wearing could save about 70,000 lives. See: <https://www.prnewswire.com/news-releases/new-ihme-covid-19-forecasts-see-nearly-300-000-deaths-by-december-1--however-consistent-mask-wearing-could-save-about-70-000-lives-301107858.html> Published on: August 6, 2020

¹⁹Tapp T (2020). Donald Trump asks reporter to take his mask off at news conference after he asks about Atlantic article. See: <https://deadline.com/2020/09/donald-trump-asks-reporter-to-take-his-mask-off-at-news-conference-atlantic-story-1234572311/> Published on: September 07, 2020

²⁰Niedzwiadek N, Atterbury A (2020). Colleges crack down on student behavior as virus threatens more closures. See: <https://www.politico.com/news/2020/08/30/college-students-coronavirus-closures-404567> Published on: August 30, 2020

²¹Moreno JE (2020). Sturgis motorcycle rally was 'superspreading event' that cost public health \$12.2 billion: analysis. See: <https://thehill.com/homenews/state-watch/515453-sturgis-motorcycle-rally-was-superspreading-event-that-cost-public> Published on: September 8, 2020

²²Groves S, Kolpack D (2020). Dakotas lead US in virus growth as both reject mask rules. See: <https://www.yahoo.com/news/virus-rises-dakotas-freedom-argument-144101521.html> Accessed on: 9/12/2020

²³The New York Times (2020). Coronavirus (Covid-19) data in the United States. See: <https://github.com/nytimes/covid-19-data> Accessed on: September 22, 2020

²⁴Raifman J, Nocka K, Jones D, Bor J, Lipson S, et al. (2020). COVID-19 US state policy database. See: www.tinyurl.com/statepolicies Accessed on: July 1, 2020

²⁵United States Census Bureau Population Division (2019). Annual Estimates of the Resident Population for the United States, Regions, States, and Puerto Rico: April 1, 2010 to July 1, 2019. 2010-2019 Population Estimates. See: <https://www.census.gov/data/tables/time-series/demo/popest/2010s-state-total.html> Accessed on: December 30, 2019

²⁶Bureau of Economic Activity (2019). Gross Domestic Product by State, 4th Quarter and Annual 2019 (PDF). See: www.bea.gov Accessed on: April 7, 2020

²⁷Kaiser Family Foundation (2018). Population Distribution by Age. See: <https://www.kff.org/other/state-indicator/distribution-by-age/?currentTimeframe=0&sortModel=%7B%22colId%22:%22Location%22,%22sort%22:%22asc%22%7D> Accessed on: December 04, 2019

²⁸Bureau of Labor Statistics (2020). State Employment and Unemployment -- JANUARY 2020. See: https://www.bls.gov/news.release/archives/laus_03162020.htm Accessed on: August 29, 2020

²⁹Department of Housing and Urban Development (2020). 2007 - 2019 PIT Counts by State (XLSX). See: <https://www.hudexchange.info/resource/3031/pit-and-hic-data-since-2007/> Accessed on: January 2020

³⁰Department of Housing and Urban Development (2020). 2007 - 2019 HIC Data by State (XLSX). See: <https://www.hudexchange.info/resource/3031/pit-and-hic-data-since-2007/> Accessed on: January 2020

³¹Kaeble D, Cowhig M and Bureau of Justice Statistics (2016). Correctional Populations In The United States, 2016. See: <https://www.bjs.gov/index.cfm?ty=pbdetail> Accessed on: August 29, 2020

³²Kaiser Family Foundation (2020). Total Number of Residents in Certified Nursing Facilities. See: <https://www.kff.org/other/state-indicator/number-of-nursing-facility-residents/?currentTimeframe=0> Accessed on: August 29, 2020

³³United States Census Bureau (2020). National Population Totals and Components of Change: 2010-2019. See: https://www.census.gov/content/census/en/data/tables/time-series/demo/popest/2010s-national-total.html#par_textimage_2011805803 Accessed on: September 23, 2020

³⁴Bennefield R, Bonnette R (2003). Structural and Occupancy Characteristics of Housing: 2000. Census 2000 Brief. See: <https://www.census.gov/prod/2003pubs/c2kbr-32.pdf> Accessed on: January 8, 2019

³⁵Kaiser Family Foundation (2020). Population distribution by race/ethnicity. See: <https://www.kff.org/other/state-indicator/distribution-by-raceethnicity/?currentTimeframe=0> Accessed on: August 29, 2020

³⁶US Congress Joint Economic Committee (JEC) Committee (2018). The Geography of Social Capital in America. See: <https://www.jec.senate.gov/public/index.cfm/republicans/2018/4/the-geography-of-social-capital-in-america> Accessed on: August 29, 2020

³⁷Centers for Disease Control and Prevention (CDC) (2018). CDC 2018 obesity map, Adult Obesity Prevalence Maps. See: <https://www.cdc.gov/obesity/data/prevalence-maps.html> Accessed on: August 29, 2020

³⁸270 to Win (2020). 2020 Gubernatorial Elections Map. See: <https://www.270towin.com/2020-governor-election/> Accessed on: September 05, 2020

- ³⁹Center for American Women in Politics (CAWP) (2020). Women in Statewide Elective Executive Office 2020. See: <https://cawp.rutgers.edu/women-statewide-elective-executive-office-2020> Accessed on: September 05, 2020
- ⁴⁰Yarnold PR and Soltysik RC (1991). Theoretical distributions of optima for univariate discrimination of random data. *Decision Sciences*, 22, 739-752.
- ⁴¹Yarnold PR, Soltysik RC (2005). Optimal data analysis: Guidebook with software for Windows. Washington, D.C.: APA Books.
- ⁴²Yarnold PR (2017). What is optimal data analysis? *Optimal Data Analysis*, 6, 26-42.
- ⁴³Soltysik RC, Yarnold PR (1994). Univariable optimal discriminant analysis: One-tailed hypotheses. *Educational and Psychological Measurement*, 54, 646-653. doi: 10.1177/0013164494054003007
- ⁴⁴Carmony L, Yarnold PR, Naeymi-Rad F (1998). One-tailed Type I error rates for balanced two-category UniODA with a random ordered attribute. *Annals of Operations Research*, 74, 223-238.
- ⁴⁵Rhodes NJ, Yarnold PR (2020). Generating novometric confidence intervals in R: Bootstrap analyses to compare model and chance ESS. *Optimal Data Analysis*, 9, 172-177.
- ⁴⁶Yarnold PR, Soltysik RC (2016). Maximizing predictive accuracy. ODA Books. doi: 10.13140/RG.2.1.1368.3286
- ⁴⁷Rhodes NJ (2020). ODA: A package and R-interface for the MegaODA software suite. R package version 1.0.1.3.
- ⁴⁸Wickham H (2016). ggplot2: Elegant graphics for data analysis. New York: Springer-Verlag. doi: 10.1007/978-3-319-24277-4
- ⁴⁹Rhodes NJ (2020). Statistical power analysis in ODA, CTA and Novometrics. *Optimal Data Analysis*, 9, 21-25.
- ⁵⁰Hendrix M (2018). The surprising geography of social capital in America. See: https://medium.com/@michael_hendrix/the-surprising-geography-of-social-capital-in-america-c6cef6bae50b Published on: June 29, 2018
- ⁵¹Yarnold PR (2016). Characterizing and circumventing Simpson's Paradox for ordered bivariate data. *Educational and Psychological Measurement*, 56, 430-442. doi: 10.1177/0013164496056003005
- ⁵²Anonymous (2020). Covid-19: New fear grips Europe as cases top 30m worldwide. See: <https://www.bbc.com/news/world-54199825> Published on: 9/18/2020
- ⁵³Scott D (2020). These 4 Midwestern states are seeing worrying Covid-19 spikes. See: <https://www.vox.com/2020/9/2/21418812/covid-19-coronavirus-us-cases-midwest-surge> Published on: September 2, 2020
- ⁵⁴Yarnold PR (2013). Ascertaining an individual patient's symptom dominance hierarchy: Analysis of raw longitudinal data induces Simpson's Paradox. *Optimal Data Analysis*, 2, 159-171.
- ⁵⁵Kretzschmar ME, Rozhnova G, Bootsma MCJ, van Boven M, van de Wijnert JHHM, et al. (2020). Impact of delays on effectiveness of contact tracing strategies for COVID-19: A modelling study. *Lancet Public Health*, 5, e452-e459. doi: 10.1016/s2468-2667(20)30157-2

⁵⁶Holcombe M, Yan H, Waldrop T (2020). ‘Mix of science and politics’ leading to people's uncertainty about Covid-19 vaccine, NIH director says. See:
<https://www.cnn.com/2020/09/16/health/us-coronavirus-wednesday/index.html> Published on: September 17, 2020

⁵⁷Makridis C, Rothwell JT (June 29, 2020). The real cost of political polarization: Evidence from the COVID-19 pandemic. doi:
10.2139/ssrn.3638373

⁵⁸<https://www.cdc.gov/mmwr/volumes/69/wr/pdfs/mm6947e2-H.pdf>

⁵⁹Tversky A, Kahneman D (1974). Judgment under uncertainty: Heuristics and biases. *Science*, 185, 1124-1131. doi:
10.1126/science.185.4157.1124

Author Notes

This study analyzed publically available data and thus was exempt from Institutional Review.