# An Observational Study to COVID-19 and 2020 US Presidential Election: Counties with High Deaths per Case Rate Showed Diminished Support For Trump*

Yiliu Cao

December 15, 2023

This study investigates the causal inference between COVID-19 and Donald Trump's loss during the 2020 US Federal Election using the data from MIT Election Data Science Lab and Johns Hopkins University CSSE. The main methodology used in this paper is propensity score matching with an exploratory analysis of the optimal treatment. The key finding suggests that the counties with a death per case rate exceeding the 0.4 quantile threshold show a reduced voting preference for Trump. This paper also conducts a counterfactual analysis based on the treatment effect, indicating that Trump might have been secured to re-elect if the disparities in death per case rates were addressed. Future research should refine the propensity score calculations by incorporating additional variables, such as the winning party in each county during the 2016 election.

## Table of contents

---

# 1 Introduction

In early 2020, the United States encountered its first cases of COVID-19. The situation escalated rapidly, with reported cases reaching 10,000 by March 19 and soaring to 100,000 just eight days later (Wikimedia Foundation 2023a). Amidst this unfolding crisis, the U.S. Presidential Election was underway, ending with Trump's defeat to Biden. By Election Day on November 3, there were more than 9 million reported cases and about 200,000 deaths in the U.S. When talking about Trump's loss, the public refers to his mishandling of "his greatest test" (Greenblatt 2021) and believes that his win in the 2016 Election was a historical accident (Bryant 2020). Conversely, some arguments suggest that COVID-19 was not the sole or decisive factor in Trump's defeat; his limits as a political strategy and wasting his advantage on the economy also contributed to his loss. This paper will investigate the causal effect between COVID and Trump's loss, and if so, to what extent it influenced the 2020 election outcome.

Existing research on this subject provides varied perspectives. Baccini et al. (Baccini, Brodeur, and Weymouth 2021) suggested that Trump could win without COVID, but simply using the COVID-19 infection rate as the only factor contributing to his loss would be too naive. Socio-economic factors like high education attainment and race diversity are also unignorable. Besides, their research emphasized the significant impact of COVID-19 on those urban areas where there were no stay-at-home orders, particularly for the "swing" states. In addition, Noland & Zhang (Noland and Zhang 2021) suggested that the deaths per case is a more crucial metric than the infection rate when analyzing the impact of COVID-19 on voting for Trump. However, they highlight the challenge in assessing COVID-19's impact on voting, as we may not know when the voters made their decisions. In contrast, Clarke et al. (Clarke,

Stewart, and Ho 2021) conducted actual surveys before and after the election, suggesting that COVID-19 impacts voting but not the dominant one. The U.S.'s highly polarized political landscape might be more significant. Similarly, Hart (Hart 2021) also conducted a survey to ask about people's attitudes toward Trump. All the participants show natural or negative possession. They also find that some social movements like "Murder of George Floyd" may impact voting.

While all the above research investigated the correlation between COVID-19 and voting for Trump, previous studies predominantly employed models to identify variables influencing Trump's vote share and perform counterfactual analyses to estimate his voting performance in the absence of COVID-19 or reduction in deaths. Instead of using that method, this study employs propensity score matching to find the treatment effect and analyze whether counties more severely affected by COVID-19 exhibit distinct voting patterns. However, given all counties have experienced COVID-19, this paper will first find the optimal treatment group(s). Finally, I will conduct a counterfactual analysis to see if Trump could be re-elected with the treatment effect eliminated.

This paper will have five parts. I will introduce the data used in this paper and present data summaries and visualization in the Section 2 part. After that, I will explain the methods in this paper and the corresponding results in the **?@sec-methods** and Section 4 parts, respectively. All the results will be interpreted and discussed in the Section 5 part, and I will conclude with limitations and drawbacks in the Section 6 part.

## 2 Data

The data in this paper is either downloaded directly or accessed via API. `R` (R Core Team 2022) will be the computer language used in the paper. Moreover, this paper will also use different R packages. The data is cleaned using the package `dplyr` (Wickham et al. 2023), `stringr` (Wickham 2022), `tidyverse` (Wickham et al. 2019), `janiotr` (Firke 2023), `tools` (R Core Team 2023) and `tidyr` (Wickham, Vaughan, and Girlich 2023). In addition to that, `ggplot2` (Wickham 2016), `RColorBrewer` (Neuwirth 2022), `maps` (Richard A. Becker, Ray Brownrigg. Enhancements by Thomas P Minka, and CRAN team. 2023), `mapdata` (Richard A. Becker and Ray Brownrigg. 2022), `gridExtra` (Auguie 2017) and `cowplot` (Wilke 2020) will be used later to make tables and plot graphs. Besides, to perform Propensity Score Matching, this paper will also employ `Matching` (Sekhon and Grieve 2012).

### 2.1 Data sources

This paper comprised five data sets from three different sources, each corresponding to different topics. The primary data is from the MIT Election Data Science Club (MIT Election Lab, n.d.), which builds open online data collections of the U.S. Federal or Senate Election

results from national to county levels. The data extracted is called "County Presidential Election Returns 2000-2020," (MIT Election Data and Science Lab 2022) with about 70,000 rows containing the voting patterns for each candidate and party by county since 2000. Besides that, the data also indicate the types of voting, such as "EARLY VOTE" and "ELECTION DAY." To analyze the voting patterns for Donald J. Trump, I only filter the 2020 U.S. Federal Election data for all counties and parties.

Additionally, the data on COVID-19 is taken from the Center for Systems Science and Engineering (CSSE) at Johns Hopkins University (CSSE 2023). CSSE collects and reports local, national, and global multidimensional data, including medicine, health care, disaster response, etc. During the pandemic, they collected the U.S. and international COVID cases and deaths, reporting them by daily summary on their GitHub (CSSEGISandData, n.d.). To examine the impact of COVID-19 on the 2020 Election as accurately as possible, this paper used the daily report on November 3, 2020, which is the election day of the 2020 U.S. Federal Election. The resulting data include the aggregate number of cases, deaths and recovers before Election Day and the incidence rate and case-fatality ratio for each county.

This paper also incorporates the socioeconomic data from the American Community Survey (ACS) (Bureau 2023), which is an online open-source database conducted by the U.S. Census Bureau, containing the various socioeconomic factors at different geographic levels. The data sets extracted from ACS are 2020 five-year estimates of DP02, DP03 and DP05, covering the social, economic and demographic characteristics at the county level. Since there are thousands of variables, to ensure simplicity, I will only use the variables that are found to be significant in predicting the COVID-19 mortality rate from my previous paper (Cao, n.d.). The descriptions of all variables can be found on Table 1.

After having five data sets, they are merged into one big data by county. Using the merged data, I calculate the percentage of votes for each party (candidate) in the 2020 election. In addition, I also create new dummy variables indicating the winning party in both elections for each county. Moreover, to accurately compare the COVID cases and deaths across all counties, I transform the number of cases and deaths to infection and mortality rate by 10,000 citizens in each county. Using the conclusion from Noland & Zhang's work (Noland and Zhang 2021), I also calculate each county's deaths per case rate. The final data consists of 3107 rows with 36 columns. All the essential variables are described in Table 1.

Table 1: Descriptions of all important variables in the analyze data

| Variables | Source | Descriptions |
| --- | --- | --- |
| income_pctile | ACS | The mean household income percentile of each county |
| prop_high_education | ACS | The proportion of local residents having a at least bachelor degree |
| private_insurance | ACS | The proportion of local residences having private insurance |
| no_insurance | ACS | The proportion of local residences without any health insurance |
| white_pct | ACS | The proportion of White population |

| Variables | Source | Descriptions |
|---|---|---|
| black_pct | ACS | The proportion of Black population |
| males | ACS | The proportion of Males |
| infrate | JHU | The COVID infection rate, calculated by cases per 10,000 residences |
| mortrate | JHU | The COVID mortality rate, calculated by deaths per 10,000 residences |
| dpc | JHU | The COVID death per 10,000 confirmed cases |
| pct_vote_demo | MIT | The percentage of votes for the Democrat in 2020 |
| pct_vote_rep | MIT | The percentage of votes for the Republican in 2020 |

## 2.2 Data summaries and visualizations

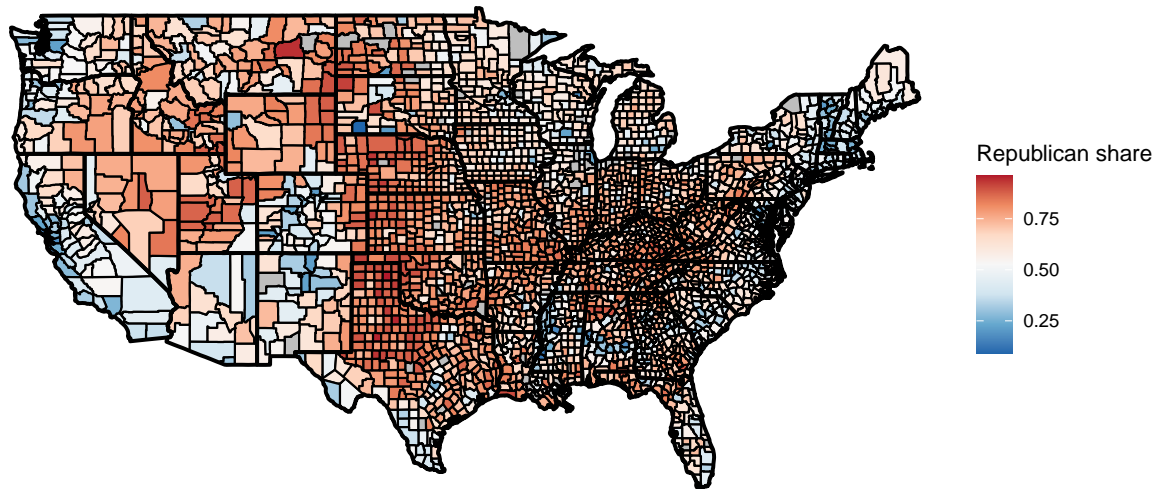Table 2: The summary of COVID cases and deaths as of Election Day.

| Winning Party | Count | Average Cases | Infection Rate | Average Deaths | Mortality Rate | DPC | Income |
|---|---|---|---|---|---|---|---|
| Democrat | 539 | 10437.514 | 2989.381 | 292.82931 | 77.394 | 2528.575 | 83540.12 |
| Republican | 2576 | 1427.246 | 2970.135 | 28.21584 | 55.643 | 1887.507 | 69800.52 |

Table 2 compares the COVID-19 impacts and income levels across the counties won by each party. Even though Republicans won in approximately five-sixths of the counties, counties that supported Biden reported nearly ten times the average number of COVID-19 cases and deaths compared to those that supported Trump. Despite similar infection and mortality rates between the two groups, there is still a notable higher death-to-case rate in counties where Democrats won. These patterns suggest that the counties voting for Democrats are more severely impacted by COVID-19 and are more population-intensive than those voting for the Republicans. Even though only around 500 counties support Democrats, they are densely populated urban areas and wealthier, so Biden can have more electoral votes. We can verify this from income levels, where the counties voting Democrat has a higher mean income.

We can validate our guess from Figure 1, which compares the relative ratio of votes for the two parties and the death to cases rate in maps. From the maps, it is clear more counties had a preference for the Republicans, but the counties with higher death-per-case rates show a preference for Democrats. These counties, such as California and New York, are mostly more affluent. In contrast, the states that are less impacted by COVID-19, such as Utah, vote more for Republicans. The maps match our findings from Table 2.

Given the insights from Table 2 and Figure 1, Figure 2 delves deeper into the impact of COVID and income levels on the voting behaviours of the two parties. The left side of Figure 2 displays the share of votes, while the right side focuses on the total number of votes, with each party

The relative share of votes between Republican and Democrat party



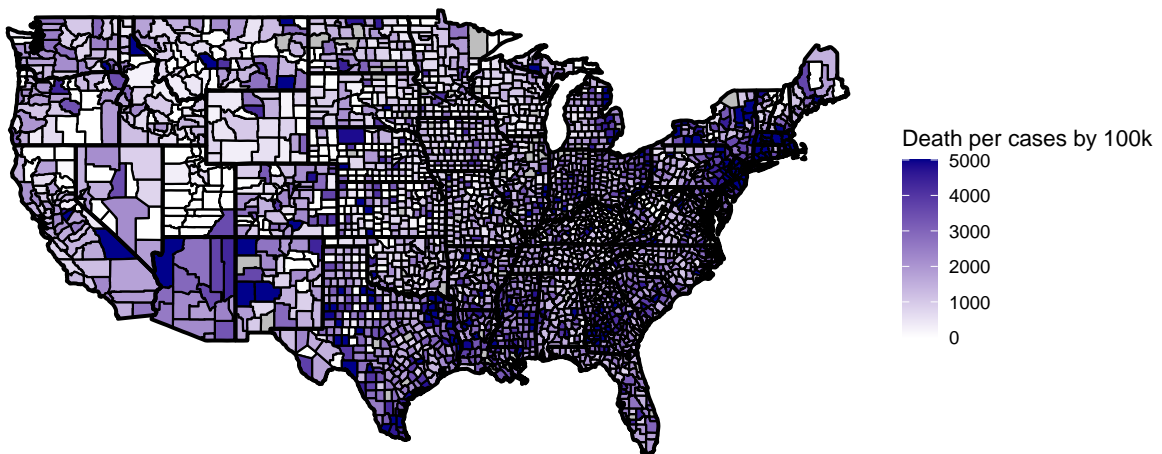The distribution of death per case rate by 100K



Figure 1: The ratio of votes for the Republican and the infection rate per 100k in each county
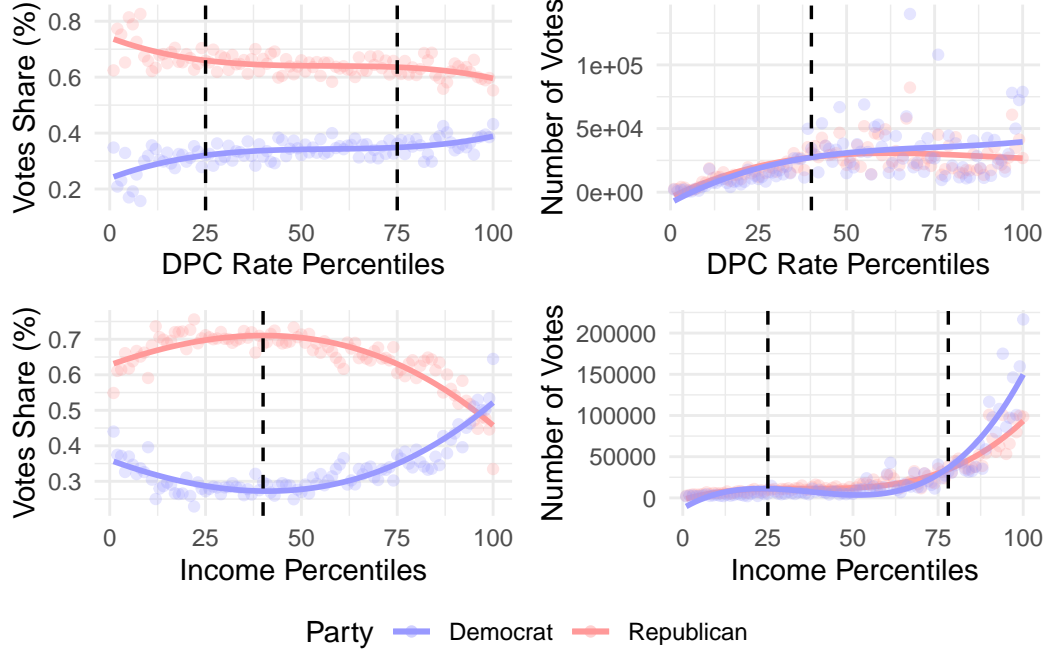
Figure 2: The correlation between voting behaviors to income and DPC rate for the Republican and Democrat

represented by distinct colours. In terms of the percentage of votes, in general, there is a positive relationship between the DPC and votes for the Democrats but a negative one for the Republicans. However, the counties with a medium value of DPC seem less sensitive to the voting for the two parties, whereas the support of counties at the 'tails' shifts more significantly with changes in the DPC. Meanwhile, there are clear quadratic patterns between the vote for the two parties and income. The lower-income counties tend to favour Republicans, but as income increases, a greater number of counties lean towards voting for Biden rather than Trump.

In addition, despite the average share of votes for Republicans being consistently above 0.5, it does not imply a universal loss for Democrats across all counties. As detailed in Table 2, we know that the number of counties won by Republicans is approximately five times greater than those won by Democrats. With Republicans dominating more counties, they naturally have a higher average share of votes. This pattern also indicates that Democrats usually won the densely populated, urban areas, but Republicans won in more sparsely populated and rural regions (Spencer 2023). The high population density areas have more votes, which explains why Biden beat Trump. We can verify this from the right side where the number of votes for Biden exceeds Trump as DPC and income level increases. This pattern is particularly pronounced in the wealthiest top 50% of counties and those with DPCs exceeding the 75th percentile.
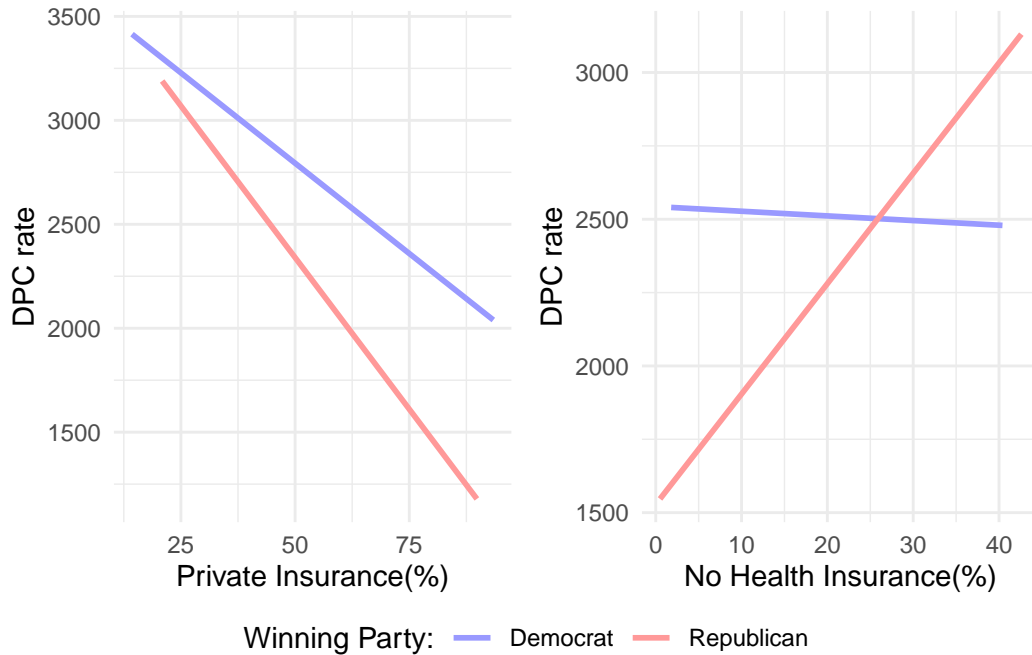
Figure 3: The correlation between DPC rate with proportons of people with private insurance and without health insurance for the two parties

Figure 3 shows the correlation between the DPC rate and the proportion of people with private and without health insurance for the two parties. Unsurprisingly, the graph reveals an inverse correlation between the DPC rate and private health insurance – as private insurance coverage increases, the DPC rate tends to decrease. Notably, counties that favoured the Democratic Party exhibit a higher proportion of individuals with private insurance compared to Republican-leaning counties. In contrast, in Republican-dominated counties, there is a pronounced positive correlation between the DPC rate and the lack of health insurance. However, it is relatively flat for Democratic counties. These findings suggest that counties voting Democratic generally have more comprehensive insurance coverage, which appears to mitigate the impact of DPC rates to a greater extent than in Republican counties.

Figure 4 compares the ratio of the White population (race diversity) and education levels for the two parties. Interestingly, it suggests that counties with higher racial diversity are less likely to vote for Biden, as opposed to the trend observed for Republican-leaning counties, where increased diversity aligns with increased support. This pattern implies that counties that predominantly vote for Democrats have a higher proportion of White residents. This may explain why they have relatively higher income levels. In terms of the ratio of people with high education attainment, since the average vote for Republicans is always higher than for Democrats, this may indicate that more educated people are likely to vote for Trump.
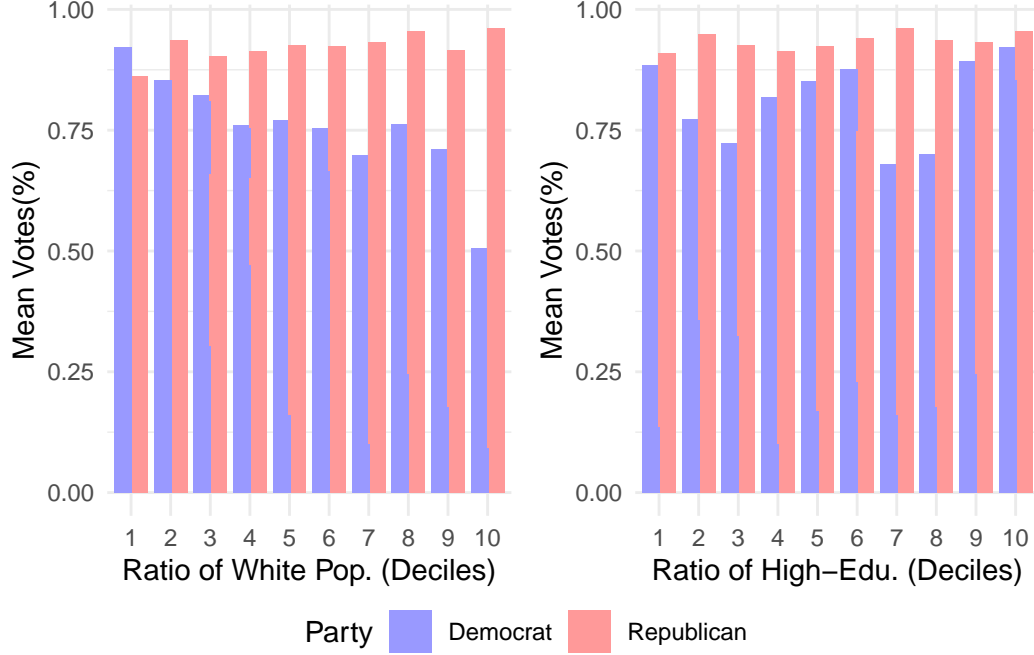
Figure 4: The distribution of share of White population and high-education across different levels

## 3 Methods

This paper aims to conduct the causal inference between COVID-19 and Trump's loss in 2020. The primary method implemented in this study will be the Propensity Score Matching (PSM). After finding the causal inference, this paper will also conduct a counterfactual analysis to see whether Trump can re-elect without COVID-19.

### 3.1 Treatment

The treatment refers to the intervention or exposure being studied to understand its potential impact on the outcome (Wikimedia Foundation 2023b). In other words, we implement specific interventions for a group of people but not for the rest. The group of people receiving the treatment is called the treatment group, and the control group for the rest. For instance, Austin (Austin 2011) examined the effectiveness of a new medical treatment called clampless off-pump coronary artery bypass (clampless OPCAB) to see if it can reduce the in-hospital mortality rate compared to the traditional method. Here, 'treatment' refers to the new surgical procedure, and patients who took this new surgery are the treatment group. They compared the mortality rates between the two groups and concluded that the new treatment could statistically lower the stroke and mortality rate.

Typically, the treatment should be defined before conducting the analysis, and this aims to reduce the risk of p-hacking (Frost 2023). However, given the objective of this paper, which is to find whether COVID-19 influenced Trump's loss in the last Election, pre-defining the treatment is challenging as all U.S. counties were affected by COVID-19. Therefore, this study begins with exploratory analysis to identify the treatment group exhibiting the most significant differences in voting patterns for Trump. I will establish specific cutoffs for classifying counties as having 'high' or 'low' death per case rates (DPC). In addition, from Figure 2, low and high-income counties also seem to have different voting behaviours. Therefore, both DPC and income will be considered when choosing the optimal treatment.

Nevertheless, choosing cutoffs is critical, and we need to avoid p-hacking or data dredging risks. P-hacking means the manipulation of data analysis until we can find statistically significant results (Frost 2023). If we simply choose one cutoff and see the significant results that match common sense, it may raise the p-hacking risk, and the analyses will be less convincing. Therefore, we will employ a grid of cutoffs for each variable, examining how treatment effects vary. For example, one treatment group can be counties with death per case and income levels higher than 0.3 quantiles. In this way, we should be able to find the optimal treatment group and avoid the risk of p-hacking.

## 3.2 Propensity Score Matching (PSM)

Randomized controlled trials (RCT) are regarded as the most ideal and the golden way to estimate the treatment effect (Kuss, Blettner, and Börgermann 2016). The reason is that RCTs are the only way that guarantee the equal distribution of known and unknown parameters. Therefore, the outcomes between the treatment and control groups will not be confounded. However, conducting an RCT is nearly impractical and impossible in many scenarios. Back to the coronary artery bypass example (Austin 2011), it's not ethically feasible to randomly assign patients to new or traditional treatments. Hence, conducting an RCT is nearly impossible.

In such situations, implementing Propensity Score Matching (PSM) can be a more practical way to find the treatment effects. The propensity score is the probability of an observation that receives the treatment based on the existing covariates using logistic regression. Using propensity scores can reduce the dimension (Zhao et al. 2021), meaning that we can describe an observation to a single score instead of multiple covariates. In addition, another significant advantage of PSM is design separation (Zhao et al. 2021), meaning that PSM can separate the covariates balancing and effects estimation. This is especially beneficial as we can observe from Figure 3 that the distribution of counties with relatively high DPC rates is imbalanced. Finally, we can match two observations having similar propensity scores but from treatment and control groups separately. Then, we will have many matched pairs and estimate the treatment effects. The matching process is called propensity score matching.

In this paper, the propensity scores are used to estimate how likely a county is to receive treatment using multiple linear logistic regression. The equation will look like:

$$\log \left( \frac{P(T = 1|X)}{1 - P(T = 1|X)} \right) = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p$$

where:
* $P(T = 1|X)$ is the probability of a county receiving treatment
* $X_j$ $(1 \leq j \leq p)$ are the covariates
* $\beta_0$ is the intercept * $\beta_j$ $(1 \leq j \leq p)$ is the corresponding coefficients of each covariate

### 3.3 Counterfactual Analysis

Using the PSM described previously, we can find the treatment effects under different treatment group settings. However, to conclude whether Trump lost due to COVID-19, we need to perform a counterfactual analysis by re-calculating the votes for Trump in counties in the treatment group.

The counterfactual analysis is particularly necessary for those "swing" states in 2020 (Wikimedia Foundation 2023c). In the 2016 Federal Election, one key factor in Trump's defeat of Hillary is that he won seven out of eleven swing states. However, he only won three in 2020. Besides, if Donald Trump were able to hold onto three swing states which are Georgia, Arizona and Wisconsin, where the average margin of Democrats is only about 0.3%, the result would have been a 269-269 electoral tie decided in the House of Representatives. The presidential election is left up to members of the House of Representatives, and Trump could win.

To find the voting patterns for Trump, especially for the "swing" states, I will follow the official election procedure and use the winners-takes-all rule (Wikimedia Foundation 2023d). That said, I will re-calculate the votes for Trump at the county level and summarize by state level. Then, the party with the highest total votes will take all the electoral votes in this state. Eventually, I will calculate the total electoral votes for each party to see whether Trump could re-elect if the treatment effects diminish.

## 4 Results

### 4.1 Choise of treatment groups

Figure 5 shows the correlation between votes for Trump and income levels, segmented into "High-DPC" and "Low-DPC" county groups at varying DPC cutoffs. For instance, a cutoff of 0.3 classifies counties with a DPC rate below this threshold as "Low DPC" and those above as "High DPC." By comparing these groups across different cutoffs, we aim to identify the cutoff where the difference in voting for Trump across different income levels is most pronounced.
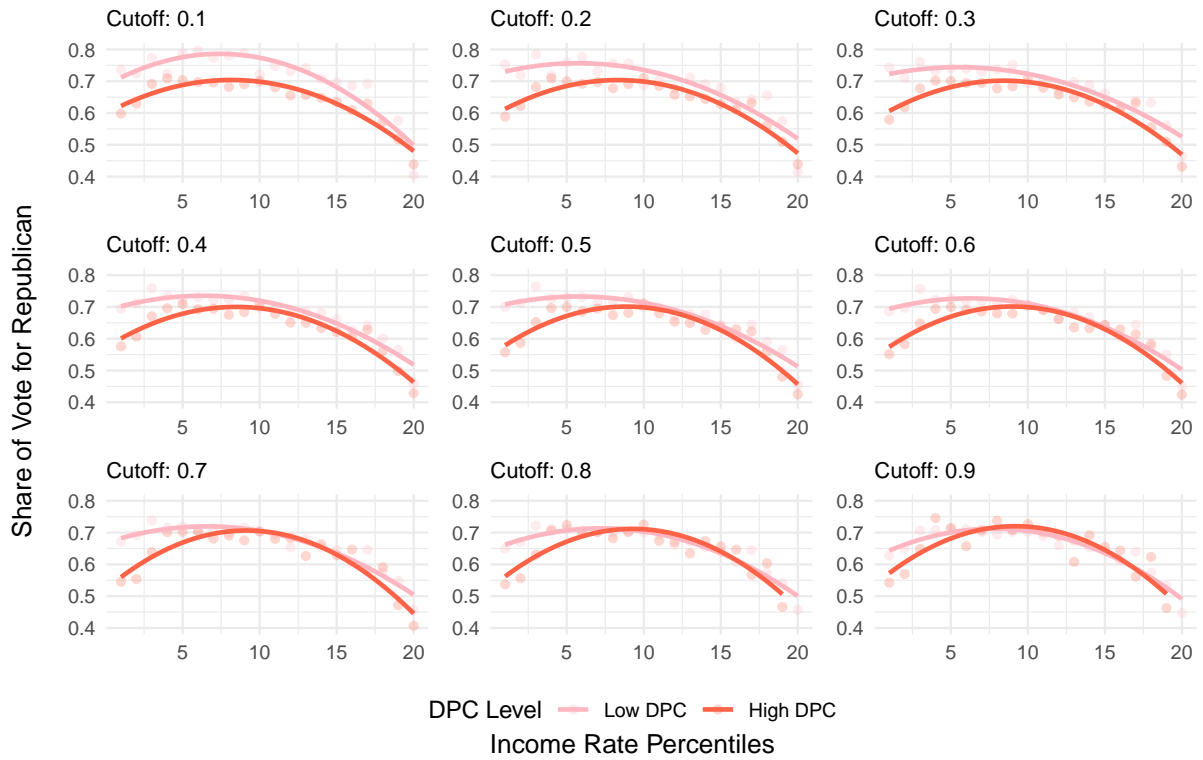
Figure 5: The Correlation Between Voting for Trump and Income Levels, Categorized by DPC Rate at Varied Quantile Cutoffs (e.g., Top 40% as High-DPC and bottom 60% as Low-DPC for "Cutoff: 0.6")

That said, we can find that, except for the cutoff 0.4, the two lines converge as income increases (right tail). Only for the 0.4 cutoff, the lines of the two groups converge at the middle income but do not overlap and diverge at the tails. Additionally, regardless of the cutoff applied, there is a common pattern that the support for Trump increases in poorer counties up to around the eighth bin, which is approximately the 0.4 income quantile, and then decreases in wealthier counties. This indicates a significant voting disparity between counties below and above the 0.4 income quantile.
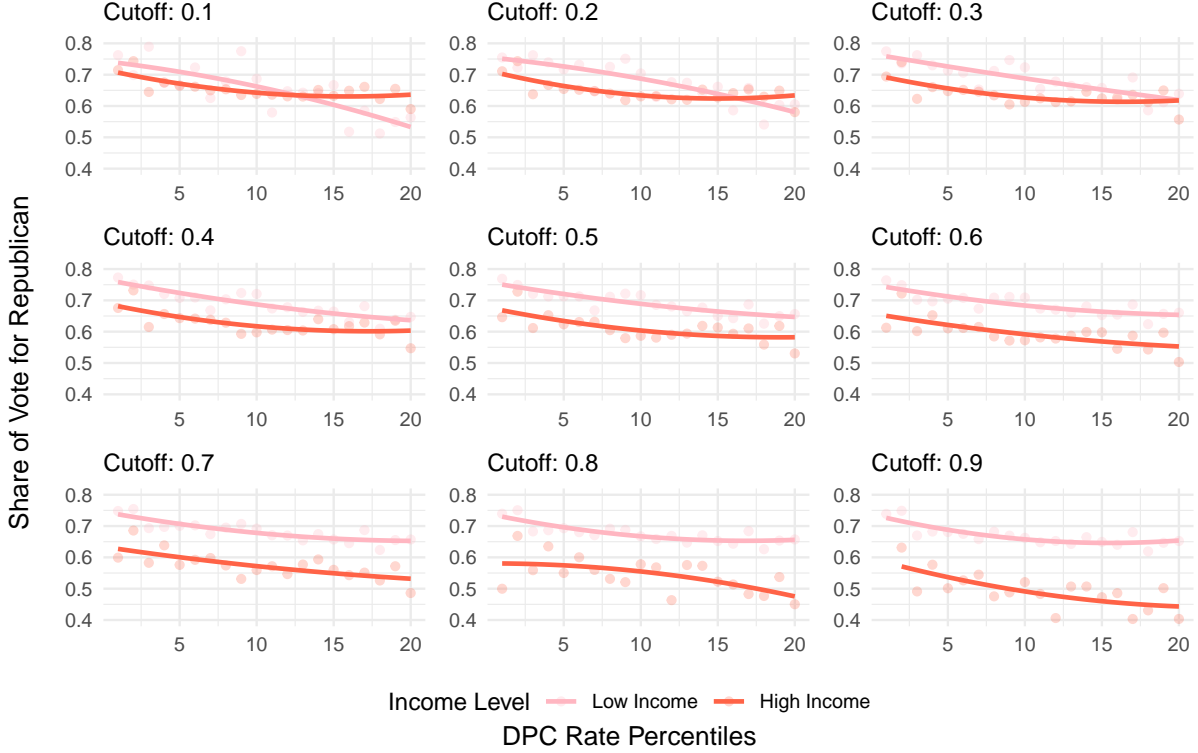


Figure 6: The Correlation Between Voting for Trump and DPC Rate, Categorized by Income Levels at Varied Quantile Cutoffs (e.g., Top 40% as High Income and bottom 60% as Low Income for "Cutoff: 0.6")

Figure 6 shows the correlation between Trump's vote share and the DPC rate but this time dividing counties based on income levels into "High Income" and "Low Income" groups. Each subplot in this analysis corresponds to a different income cutoff. Compared to Figure 5, all curves, regardless of the cutoff, display a linear or near-linear relationship with voting patterns. As the income cutoff increases, the difference in voting behaviour between the high and low-income counties becomes more distinct, especially at the 0.9 cutoff. This suggests that in counties with higher incomes, voting patterns vary significantly from those in lower-income counties as the DPC rate changes. Therefore, the most significant pattern is observed at the 0.9 income cutoff, contradicting the patterns identified in Figure 5.

So, a question may arise regarding whether we should include all the income and DPC in the settings of treatment. The answer is NO, and we should only include the DPC. From Figure 5, a consistent quadratic relationship between income and voting patterns for Trump is evident regardless of how the cutoff for defining high DPC is adjusted. This pattern demonstrates that poorer counties (in the first half of the income spectrum) may show a positive or a negative correlation with voting for Trump depending on the cutoff we choose. In contrast, more affluent counties will only exhibit a negative correlation no matter which cutoff we choose. This quadratic pattern can be verified from Figure 2, which illustrates a similar relationship between votes and income.

Conversely, Figure 6 presents a different narrative. All regression lines, irrespective of the DPC cutoff, display a linear or nearly linear relationship. We should exclude income because the choice of cutoff will significantly change the treatment effects; DPC is more "stable" than income. We can verify this from Figure 6, where each line shows a negative slope, meaning that the choice of cutoff for the "High DPC" group will not change the general voting patterns for Trump. Hence, including income would introduce extra variability in the treatment effects and increase the probability of p-hacking. The Figure 7 shows a simple visualization of why this is the case.
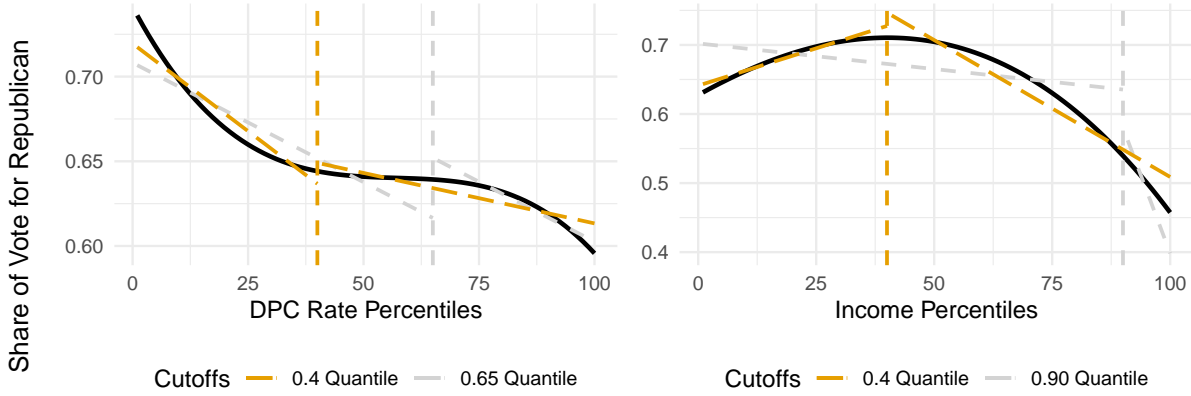


Figure 7: The simple visualization regarding why excluding income from treatment

Figure 7 shows a simple example when we test different cutoffs for DPC rate and income. From the graph, we can see that no matter which cutoff we set, the low DPC and high DPC counties always have a negative value of slope. The only difference is the absolute values of the slope. Conversely, the impact of income levels shows notable variability. For instance, at a cutoff of the 0.4 quantiles, low and high-income groups display opposite voting patterns for Trump. However, when we increase the cutoff to the 0.9 quantile, both income groups negatively correlate with Trump's vote share. This underscores the instability of income.

From Figure 5, we have determined that a cutoff at the 0.4 quantile is the most effective threshold for distinguishing between high and low Death Per Case (DPC) rate groups. Therefore, in our study, the treatment group is defined as those counties where the DPC rate exceeds the 0.4 quantile. Moreover, while income is not considered for determining the treatment and

Table 3: The summary of Propensity Score Matching

|  |  | Original | | Matched | |
| --- | --- | --- | --- | --- | --- |
| Treatment Effect | P value | Obs. | Treat. Obs. | Treat. Obs | Treat. Obs. Unw. |
| -0.018342 | 0.021404 | 3115 | 1869 | 1869 | 5415 |

control groups due to its variability, its influence is still significant and shows an imbalance from Figure 5. To account for this, income will be integrated into the calculation of propensity scores for each county.

## 4.2 Propensity Score Matching

In my previous paper, I already showed that the proportion of residents with at least a bachelor's degree, income, the proportion of people with private insurance and without health insurance, the ratio of males, ratio of black and black residents are statistically in predicting the probability of a county has a high mortality rate or not. In this paper, I will directly use the results and only consider these variables when predicting the propensity scores.

$$
\log\left(\frac{P(T = 1|X)}{1 - P(T = 1|X)}\right) = 8.06
$$

$$
- 0.02 \times \mathrm{prop\_higher\_education}
$$
$$
+ 0.01 \times \mathrm{pctile}
$$
$$
+ 0.01 \times \mathrm{no\_insurance}
$$
$$
- 0.03 \times \mathrm{private\_insurance}
$$
$$
- 0.13 \times \mathrm{males}
$$
$$
+ 0.00 \times \mathrm{white\_pct}
$$
$$
+ 0.06 \times \mathrm{black\_pct}
$$

The above equation shows the logistic regression model to predict the propensity score for each county. Males seem to be the most critical factor when predicting the propensity score. Perhaps the majority of COVID deaths are males. Interestingly, the counties with a higher income level will have a higher death per case rate. This may explain why richer counties are less likely to vote for Trump.

Table 3 summarize the results of Propensity Score Matching. The treatment effect of about -0.018 means that the counties with a DPC rate higher than 0.4 will, on average, vote 0.018 less for Trump compared to the counties with a lower DPC rate. This value is statistically significant as its p-value is about 0.02, lower than 0.05. In addition, there are 1869 treatment

Table 4: The balance of each covariates before and after the mathcing

| | Before Matching | | After Matching | |
|---|---|---|---|---|
| Variables | Mean Contr | Std Mean Diff | Mean Contr. | Std Mean Diff. |
| prop_higher_education | 23.2460 | -10.772 | 21.6330 | 5.6159 |
| pctile | 51.7740 | -10.389 | 47.8040 | 2.9221 |
| no_insurance | 8.7107 | 24.449 | 9.7494 | 4.6766 |
| private_insurance | 67.6560 | -35.913 | 63.4090 | 5.9550 |
| males | 50.5120 | -32.492 | 49.5520 | 9.6593 |
| white_pct | 87.1370 | -47.207 | 79.6360 | -5.5169 |
| black_pct | 4.1685 | 49.083 | 11.4560 | 5.4646 |

observations in both original and matched data, meaning that some counties in the control group are matched more than once.

Table 4 summarizes each variable's mean values in the control group and standard mean difference before and after the matching. Notably, there is a significant decrease in the standard mean difference for each variable. This indicates that the Propensity Score Matching effectively reduced the imbalance between the treatment and control groups. Besides, we can see that the mean values in the control group approximately remain the same after the matching, enhancing the reliability of the matching process and suggesting that PSM has successfully aligned the groups based on the observed covariates.

## 4.3 Counterfactual Analysis

Table 5 provides a comparative summary of the actual and counterfactual voting results in eleven swing states during the 2020 U.S. Presidential election. In reality, Trump only won three swing states with 232 electoral votes. However, the counterfactual analysis suggests that Trump would have won five additional states if the DPC rate disparities were eliminated, bringing him an extra 52 electoral votes. This would bring his total to 284 electoral votes, surpassing the required 270-vote threshold for reelection. Therefore, if Trump can eliminate the extra deaths, he could be re-elected.

16

Table 5: The summary of counterfactual votes for the eleven swing staes in 2020

| States | Electoral Votes | Actual Results | | | Simulated Results* | | |
|---|---|---|---|---|---|---|---|
| | | Trump | Biden | Margin | Trump* | Biden* | Margin* |
| NH | 4 | 365660 | 424937 | 7.5% D | 378054.6 | 412542.4 | 4.36% D |
| MN | 10 | 1484065 | 1717077 | 7.28% D | 1514486.2 | 1686655.8 | 5.38% D |
| MI | 15 | 2649852 | 2804040 | 2.83% D | 2733394.7 | 2720497.3 | 0.24% R |
| NV | 6 | 669608 | 703314 | 2.46% D | 693512.9 | 679409.1 | 1.03% R |
| PA | 19 | 3377674 | 3458229 | 1.18% D | 3497988.6 | 3337914.4 | 2.34% R |
| WI | 10 | 1610065 | 1630673 | 0.64% D | 1613541.4 | 1627196.6 | 0.42% D |
| AZ | 11 | 1661686 | 1672143 | 0.31% D | 1723779.1 | 1610049.9 | 3.41% R |
| GA | 16 | 2461837 | 2474507 | 0.26% D | 2547370.9 | 2388973.1 | 3.21% R |
| NC | 16 | 2758773 | 2684292 | 1.37% R | 2810662.1 | 2632402.9 | 3.27% R |
| FL | 30 | 5668731 | 5297045 | 3.39% R | 5847431.0 | 5118345.0 | 6.65% R |
| IA | 6 | 902009 | 759061 | 8.61% R | 917420.3 | 743649.7 | 10.46% R |

# 5  Discussion

## 5.1  First discussion point

## 5.2  Second discussion point

## 5.3  Third discussion point

## 5.4  Weaknesses and next steps

# 6  Conclusion

# References

Auguie, Baptiste. 2017. *gridExtra: Miscellaneous Functions for "Grid" Graphics.* https://CRAN.R-project.org/package=gridExtra.

Austin, P. C. 2011. "An Introduction to Propensity Score Methods for Reducing the Effects of Confounding in Observational Studies." *Multivariate Behavioral Research*, May. https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3144483/.

Baccini, L., A. Brodeur, and S. Weymouth. 2021. "The COVID-19 Pandemic and the 2020 US Presidential Election." *Journal of Population Economics.* https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7809554/.

Bryant, N. 2020. "US Election 2020: Why Donald Trump Lost." https://www.bbc.com/news/election-us-2020-54788636.

Bureau, U. S. Census. 2023. "American Community Survey (ACS)." 2023. https://www.census.gov/programs-surveys/acs.

Cao, Y. n.d. "Yiliuc/Covid_us_county: This Repository Contains the Paper of COVID-19 Mortality Rate in STA497." https://github.com/yiliuc/COVID_US_county/tree/main.

Clarke, H., M. C. Stewart, and K. Ho. 2021. "Did Covid-19 Kill Trump Politically? The Pandemic and Voting in the 2020 Presidential Election." *Social Science Quarterly.* https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8242570/.

CSSE, JHU. 2023. "COVID-19 Content Portal." 2023. https://systems.jhu.edu/research/public-health/ncov/.

CSSEGISandData. n.d. "CSSEGISANDDATA/COVID-19: Novel Coronavirus (COVID-19) Cases, Provided by JHU CSSE." https://github.com/CSSEGISandData/COVID-19.

Firke, Sam. 2023. *Janitor: Simple Tools for Examining and Cleaning Dirty Data.* https://CRAN.R-project.org/package=janitor.

Frost, J. 2023. "What Is p Hacking: Methods & Best Practices." https://statisticsbyjim.com/hypothesis-testing/p-hacking/.

Greenblatt, A. 2021. "Five Reasons Donald Trump Lost the Presidency." https://www.governing.com/now/five-reasons-donald-trump-lost-the-presidency.html.

Hart, J. 2021. "Did the COVID-19 Pandemic Help or Hurt Donald Trump's Political Fortunes?" *PloS ONE.* https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7904340/.

Kuss, O., M. Blettner, and J. Börgermann. 2016. "Propensity Score: An Alternative Method of Analyzing Treatment Effects." *Deutsches Ärzteblatt International.* https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5963493/#R16.

MIT Election Data and Science Lab. 2022. "County Presidential Election Returns 2000-2020." Harvard Dataverse. https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/VOQCHQ.

MIT Election Lab. n.d. "Home Page." https://electionlab.mit.edu/.

Neuwirth, Erich. 2022. *RColorBrewer: ColorBrewer Palettes.* https://CRAN.R-project.org/package=RColorBrewer.

Noland, M., and Y. E. Zhang. 2021. "Covid-19 and the 2020 US Presidential Election: Did the Pandemic Cost Donald Trump Reelection?" https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3807255.

R Core Team. 2022. *R: A Language and Environment for Statistical Computing.* Vienna, Austria: R Foundation for Statistical Computing. https://www.R-project.org/.

———. 2023. *R: A Language and Environment for Statistical Computing.* Vienna, Austria: R Foundation for Statistical Computing. https://www.R-project.org/.

Richard A. Becker, Original S code by, and Allan R. Wilks. R version by Ray Brownrigg. 2022. *Mapdata: Extra Map Databases.* https://CRAN.R-project.org/package=mapdata.

Richard A. Becker, Original S code by, Allan R. Wilks. R version by Ray Brownrigg. Enhancements by Thomas P Minka, and Alex Deckmyn. Fixes by the CRAN team. 2023. *Maps: Draw Geographical Maps.* https://CRAN.R-project.org/package=maps.

Sekhon, Jasjeet Singh, and Richard D. Grieve. 2012. "A Matching Method for Improving Covariate Balance in Cost-Effectiveness Analyses." *Health Economics* 21 (6): 695–714.

Spencer, S. H. 2023. "Number of Counties Won in Presidential Election Doesn't Determine Outcome." https://www.factcheck.org/2023/08/number-of-counties-won-in-presidential-election-doesnt-determine-outcome/.

Wickham, Hadley. 2016. *Ggplot2: Elegant Graphics for Data Analysis.* Springer-Verlag New York. https://ggplot2.tidyverse.org.

———. 2022. *Stringr: Simple, Consistent Wrappers for Common String Operations.* https://CRAN.R-project.org/package=stringr.

Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D'Agostino McGowan, Romain François, Garrett Grolemund, et al. 2019. "Welcome to the tidyverse." *Journal of Open Source Software* 4 (43): 1686. https://doi.org/10.21105/joss.01686.

Wickham, Hadley, Romain François, Lionel Henry, Kirill Müller, and Davis Vaughan. 2023. *Dplyr: A Grammar of Data Manipulation.* https://CRAN.R-project.org/package=dplyr.

Wickham, Hadley, Davis Vaughan, and Maximilian Girlich. 2023. *Tidyr: Tidy Messy Data.* https://CRAN.R-project.org/package=tidyr.

Wikimedia Foundation. 2023a. "Covid-19 Pandemic in the United States." https://en.wikipedia.org/wiki/COVID-19_pandemic_in_the_United_States.

———. 2023b. "Propensity Score Matching." https://en.wikipedia.org/wiki/Propensity_score_matching.

———. 2023c. "Swing State." https://en.wikipedia.org/wiki/Swing_state#cite_ref-tp2020_42-0.

———. 2023d. "Winner-Take-All Market." https://en.wikipedia.org/wiki/Winner-take-all_market.

Wilke, Claus O. 2020. *Cowplot: Streamlined Plot Theme and Plot Annotations for 'Ggplot2'.* https://CRAN.R-project.org/package=cowplot.

Zhao, Q.-Y., J.-C. Luo, Y. Su, Y.-J. Zhang, G.-W. Tu, and Z. Luo. 2021. "Propensity Score Matching with r: Conventional Methods and New Features." *Annals of Translational Medicine*, May. https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8246231/.