Yichuan Shi

17.803

3/21/2022

Design Project Final Report

**Introduction**

Stop Enabling Sex Traffickers Act (SESTA) and Allow States and Victims to Fight Online Sex Trafficking Act (FOSTA) are laws passed in 2018 that stripped Section 230's safe harbor provisions in instances of sexual trafficking. While websites and platform providers are usually not legally liable for user-generated content, the SESTA-FOSTA bill ensured that such protections were stripped in cases of suspected sex-trafficking. Immediately following the passage of these bills, the largest classified ads page in the US, [backpage.com](backpage.com), was seized and taken down by the FBI on grounds of enabling human trafficking. Moreover, online platforms in general began to preemptively filter and remove sexually-explicit content and listings in fear of legal blowback.

While proponents of SESTA-FOSTA argued that the legislations effectively shut down one of the largest platforms for human trafficking in the US, opponents claimed the bills' broad brushstroke and vague language also stiffled online resources used by sex-workers to stay safe. Without online platforms and crowdsourced webpages to vet clientele, sex workers would resort to more traditional means of door-knocking to find work. As a result — according to activist groups — instances of violence against sex workers also skyrocked. While both sides of the debate have released qualitative studies and memos, there exist no comprehensive study on the causal impact of SESTA-FOSTA's stripping of safe harbor provisions on trafficking rate and violence against sex workers.

Our hypothesis is that SESTA-FOSTA's passage resulted in an increase of violent crimes against sex workers. The null hypothesis is that SESTA-FOSTA's passage had no impact on violent crimes against sex workers. The significance level chosen to reject the null hypothesis is chosen at $p = 0.05$.

**Theoretical Motivation**

The case's theoretical importance spans across regulatory debates concerning both sex-work and online platform providers. Firstly, understanding the causal impact of removing sex-workers' online platforms on violence against sex-workers clarifies how to best regulate sex work. Currently, the regulatory question is split between curbing supply versus curbing demand:[1] The former camp of law-makers emphasize on regulating providers of sex work, whereas the latter camp focuses on the "Johns" that consume the service. Policies targeting the demand side of sex-work include sting operations, reeducation, and shaming through "Dear John" letters sent to family members.[2] On the other hand, stripping away platform protection for explicit listings to minimize content related to sex-work — as SESTA-FOSTA had done — falls squarely into supply-targeted regulation. As a result, SESTA-FOSTA's effectiveness directly reflect whether supply-focused approaches are effective in protecting vulnerable populations against sexual exploitation.

Moreover, SESTA-FOSTA's regulatory mechanism has deep implications for internet regulation and corporate balance in the US. Section 230 of the Communications Decency Act had long been used to shield platforms from legal consequences of user-generated content, which in turn inspired features that defined the internet, such as discussion forums and social media platforms. Stripping the safe harbor provisions through explicit content, however, resulted in waves of self-censorship by smaller discussion boards and platforms.[3] While sex-workers and community-based activists and organizations overwhelming oppose SESTA-FOSTA, tech companies also briefly lobbied against the legislation due to increased legal responsibility. On the other hand, corporate giants like Disney and 20th Century Fox lobbied in support of SESTA-FOSTA, precisely because the legislation would pave the way for automatic filtering and censoring of copy-righted content.[4] In other words, SESTA-FOSTA provided important legal precedence for platform oversight, self-regulation, and market balance on the internet. Moreover, despite significant monetary interest on both sides of the debate by corporate giants, sex-workers would shoulder the bulk of the fallout of legislative impacts.

**Literature review**

There exists numerous accounts of well-established literature that detail the progression, consensus, philosophical underpinning, and debates on the regulation of human-trafficking and sex-work. Most important to this case-study, however, are research that details precise indicators of sex-work, the actual platforms that were impacted by SESTA-FOSTA, and the specific alternatives that sex-workers resorted to in the wake of SESTA-FOSTA's passage. Knowing such quantifiable information allows us to verify the internal validity of the hypothesis, and to perform data analysis that targets sex-workers as the treated group.

Firstly, the New York City Trafficking Assessment Project[5] and the United Nations Office on Drugs and Crime[6] converge on several indicators of interest in their independent findings — such as frequent change of location, living with non-family members, limited interaction with those outside of immediate living arrangements, not speaking the local language, no cash of their own, no access to personal legal documents, so on and so forth. More specifically, NYC's Trafficking Assessment Project culminated in a screening toolkit with specific questions that poke at indicators of sex-work and human-trafficking victimization. Specific indicators from the screening toolkit that correspond to fields of interest in the dataset include type of living quarter, whether the living situation is a gated or walled community, household income, marital status, educational attainment, and more.

To validate the hypothesis' internal validity, it is necessary to understand the type of platforms affected by SESTA-FOSTA, and the alternative means of finding work that sex-workers resorted to. According to an online survey performed by the Anti-Trafficking Review, seventy percent of sex-workers reported a change in online resources avaliable since the law's passing.[7] Moreover, one third of respondents reported an increase of violence from clients in the same survey. Specific platforms that removed resources include Google, Microsoft, Reddit, Skype, and Craigslist.[8] For instance, Google Docs have implemented automatic detection and deletion of NSFW content, Craigslist have removed escort listings, and numerous sex-work related subreddits have been removed due to fear of legal reprecussion.[9]

Popular resources for vetting Johns such as VerifyHim have also been removed in the wake of SESTA-FOSTA.[10]

As a result of the deplatforming, sex-workers have resorted to traditional, in-person alternatives to find work. According to the Berkeley Journal of Criminal Law, crimes related to pimping and street-based transactions more than tripled in San Francisco.[11] From interviews with police officers, it was also revealed that former pimps of emancipated sex-workers reemerged after the passage of SESTA-FOSTA with the specific intention of re-recruiting former victims. Moreover, New York City in the same year represented an almost two-fold increase in arrests for loitering for prostitution, concentrated in immigrant-heavy neighborhoods.[12] Overall, it's evident that the stripping of online resources put sex-workers in a much more dangerous alternative.

**Empirical Strategy**

To test the hypothesis, this research will analyze data from the National Crime Victimization Survey (NCVS).[13] Hosted by the Bureau of Justice Statistics, the NCVS is an annual survey conducted on a representative sample of about 240,000 people from around 150,000 households. NCVS data collection has been ongoing since 1973 and the latest available data is from 2020. Interviewees are asked about crime indicators from nonfatal personal crimes, financial crime, property crime, and violent crimes. Infromation on age, sex, race, marital status, education level, household location, and income are collected from the respondents. The specific data fields of interest for the purpose of this study are detailed in the earlier literature review section. Moreover, the NCVS asks questions on whether the crime was reported, reasons for non-reporting (if that were the case), and victim experiences with the criminal justice system.

Besides detailed questions on respondents' background and crime victimization experiences, the NCVS is also unique in how continuity is reflected in its results. A household is defined as any group of people who reside in the same address. The selection process is also weighted to adjust to known population totals and compensate for non-response. Once selected, the residence stays in the survey for

3.5 years and interviewed every 6 months — in other words, each household occupies seven data points. New households rotate into the sample set as old households leave after 7 interviews.

Overall, the NCVS survey format and extensive data allows us to subset household and interviewee by characteristics that approximate which respondents might be sex-workers, using the indicators revealed by the literature review. Moreover, the household rotation schema allows us to subset the records by data points that span before and after the enactment of SESTA-FOSTA, which was signed into law on April 11, 2018. NCVS's household rotation system allows for tracking on the treatment's specific impacts on the same respondents.

**Research Design**

The proposed research design is a difference-in-differences analysis over all households in the NCVS dataset with records that overlap the enactment of SESTA-FOSTA, matched with identified covariates. The treatment in question is SESTA-FOSTA. The treated group is the approximated population for sex-workers, estimated by indicators outlined in the literature review section. Finally, the control group is every household that does not share the same indicator. All members of the treated group received the treatment, because SESTA-FOSTA was enacted and enforced across the US. Moreover, over the NCVS's large dataset, individual household differences can be amortized and population-wide treatment effects can be gauged. If the changes in personal and violent crime rates before and after 2018 differ significantly between the treated and control groups, we can reasonably conclude the treatment's causal impact.

Difference-in-differences is chosen as the method of analysis because many other potentially confounding variables also occurred during the timeframe in question, the most significant of which being the 2019 Coronavirus Pandemic. Pandemic-era industry disruption and displacement lead to a sharp uptick in rates of homicide, aggravated assault, and firearm-related assault.[14] In fact, homicides increased by almost fifty percent compared to before the pandemic. Moreover, unlike disruptive events in the past, smaller towns and rural areas experienced the same spike in crime rates as urban centers did.[15] It is

paramount for the analysis to account for the confounding factor's major impact on our outcome variables of interest.

Preprocessing with Matching is chosen because the ratio between the control dataset and the treatment dataset is incredibly large– the cleaned dataset before matching contained 139,862 control data points and 5,395 treatment data points, for a total control-to-treatment ratio of around 26:1. The high ratio means that we can comfortably perform matching without diminishing external validity by throwing away too many data points. Moreover, unit-to-unit matching results in a more precise control-to-treatment balance that does not make strong assumptions about similarities. In fact, there are a lot of interview results in the control group that vary drastically in personal background and circumstances from the treatment group, which means that the *Ceteris Paribu*s assumption really cannot be comfortably applied. As a result, one-to-one matching with replacement allows us to account for some bias and variance in the dataset, as well as avoid making strong assumptions about the similarities between the treatment and control group. This is especially necessary because the treatment variable is imprecise in the first place.

**Analysis and Findings**

The NCVS dataset used has 297,399 entries across 1,017 data fields. The unit of analysis is a single interview with a single respondent who reports some form of crime in the time range between 1992 and 2020. The primary outcome variable of interest is the field labeled V4529, which measures the type of crime code and is present for all years and quarters except for 1992 Q1 and Q2.

*Identification of Treatment and Covariates*

The treatment variable in our analysis is whether someone is a sex-worker. Obviously, the NCVS dataset does not directly ask respondents this question, and we must approximate the respondent's treatment status based on indicator variables from the other recorded datafields. Through previously cited literature — specifically the qualitative interviews with small groups of sex-workers such as studies by the Urban Justice Center and Whose Corner Is It Anyway (WCIIA) — we can approximate a given

respondent's sex-worker status by variables like their living situation, financial status, frequency of reallocation, prior interactions with police, type of living quarters, and more. The table below documents the procedure used to determine whether a given data entry is treatment or control.

| IF (V2014 = 2 or 3) | Rented unit for cash or no cash rent | AND |
|---|---|---|
| (V2020 = 1, 2, 3, 4, 5, 6, 7, 8, 9, 10) | Type of living quarters is house, apartment, flat, non-transient hotel, motel, transient hotel, motel, rooming house, mobile home or trailer, boarding house | AND |
| (V2024 != 98, 99) | Number of units in structure is not residue or out of universe | AND |
| (V2025A != 2) | Not living in a gated/walled community | AND |
| (V2025B != 2) | Building does not have restricted access | AND |
| (16 <= V2033 <= 60) | Aged between 16 to 60 | AND |
| (V2036 = 2 or trans = TRUE) | Either female or trans-identifying | AND |
| (V2122 contains "lone") | Living alone partner in family structure code | AND |
| (V3033 > 2) | Moved more than twice in the past five years | AND |
| (V2026 <50000) | Annual household income less than $50,000 | OR |
| (V2074 = 1 and V2075 = 2) | Operate business from address but does not have a sign | THEN |
| Treat = 1 | Respondent is categorized as a sex-worker | |

Besides the categories used to identify the treatment group, three variables are also singled out as potential covariates: education, urban or rural location, and race. Since these variables are not used to identify the treatment, we can match the dataset on these covariates. A later analysis specifically looks at race's role on sexual assualt rates pre and post SESTA-FOSTA, so the dataset used for that separate analysis is not matched on race.

Before any statistical analysis or pre-processing, I graphed the treatment and control groups across the three outcome variables of interest to test the parallel trend assumption. As Figure 1 indicates, both assault and sexual assault display reasonable parallel trend, whereas robbery does not.
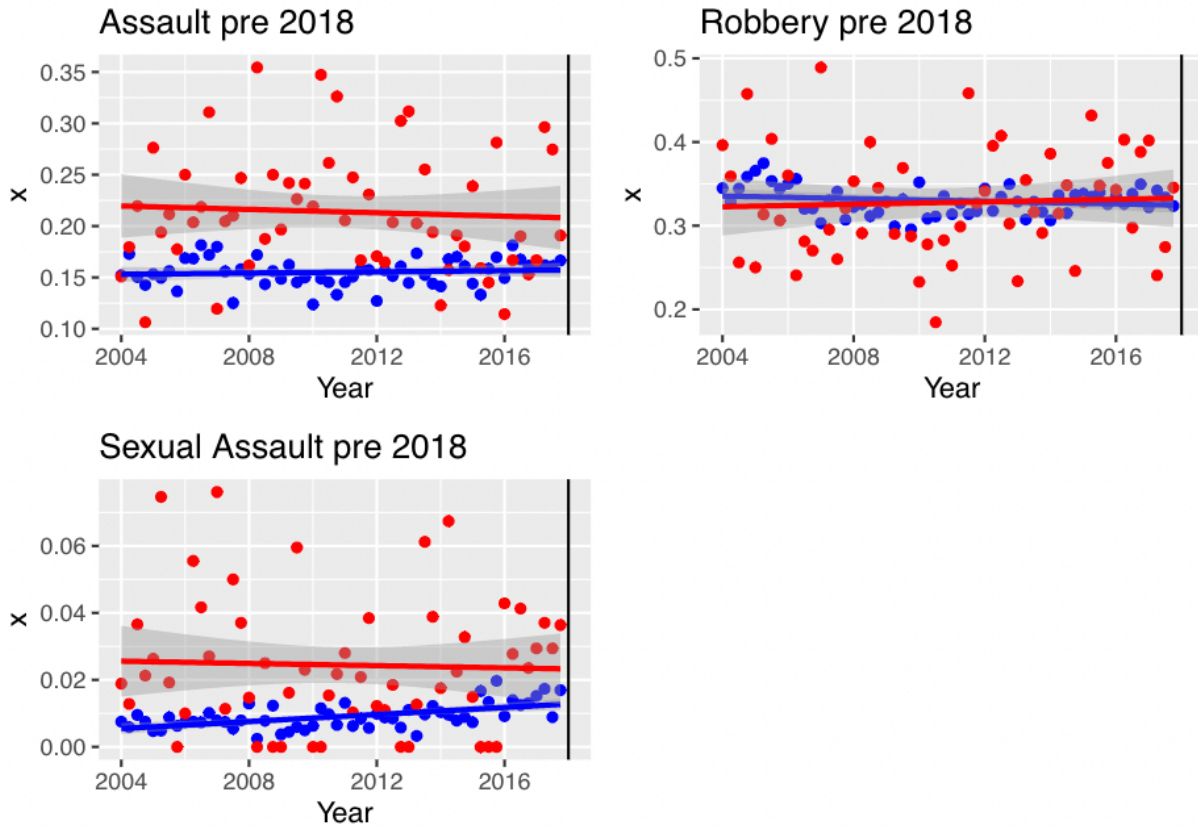
*Figure 1, Parallel trend diagnostics of data across three outcome variables of interest pre SESTA-FOSTA*

*Pre-Processing with Matching*

Table 1 shows the balance of covariates across the three outcome variables of interest before matching. Figure 2 in the appendix shows As displayed in the T pval column, all three covariates across all three outcome fields are statistically significant. This means that matching was the correct decision to take, and that without account for the presence of our identified confounders, we cannot cleanly attribute the outcome to be solely due to the treatment.

Table 1: Balance Pre-Matching

|  | mean.Tr | mean.Co | T pval |
|---|---|---|---|
| education.robbery | 20.33 | 21.29 | 0.00 |
| urban.rural.robbery | 1.10 | 1.16 | 0.00 |
| race.robbery | 1.55 | 1.44 | 0.00 |
| education.assault | 20.33 | 21.29 | 0.00 |
| urban.rural.assault | 1.10 | 1.16 | 0.00 |
| race.assault | 1.55 | 1.44 | 0.00 |
| education.sexual.assault | 20.33 | 21.29 | 0.00 |
| urban.rural.sexual.assault | 1.10 | 1.16 | 0.00 |
| race.sexual.assault | 1.55 | 1.44 | 0.00 |

*Table 1: Balance of covariates across three outcome variables pre-matching*

After matching with Mahalanobis distance with replacement and with randomly-broken ties to speed up performance, we see that the covariate balance improved significantly. Now, none of the covariates are statistically significant, and we can more confidently attribute results from the outcome to be solely due to the treatment. Moreover, Figure 3 in the appendix shows the density balance of covariates after matching, and the density overlap becomes a lot more coherent after matching. With covariates properly balanced across the treatment and control units, we can finally begin the difference-in-differences analysis!

Table 2: Balance Post-Matching

|  | mean.Tr | mean.Co | T pval |
|---|---|---|---|
| education.robbery | 20.33 | 20.33 | 0.08 |
| urban.rural.robbery | 1.10 | 1.10 | 1.00 |
| race.robbery | 1.55 | 1.55 | 0.08 |
| education.assault | 20.33 | 20.33 | 0.56 |
| urban.rural.assault | 1.10 | 1.10 | 1.00 |
| race.assault | 1.55 | 1.55 | 0.08 |
| education.sexual.assault | 20.33 | 20.33 | 0.08 |
| urban.rural.sexual.assault | 1.10 | 1.10 | 1.00 |
| race.sexual.assault | 1.55 | 1.55 | 0.08 |

*Table 2: Balance of covariates across three outcome variables post-matching*

*Difference-in-Differences results*

Three difference-in-differences tests were performed with the three matched datasets across the three outcome variables of interest. The results are summarized in table 3 — While robbery and assault outcomes do not display any statistical significance, sexaul assult outcomes are statistically significant at the 0.05 significance level. This means that we have a five percent chance of rejecting the null hypothesis incorrectly, of concluding that there exists statistically significant relationship of SESTA-FOSTA enactment and sexual assault rates among sex workers, when in fact there are none. The observation size in our analysis is $N = 1445$. We can conclude that the evidence in our sample size is strong enough to reject the null hypothesis at the population level. The implementation of SESTA-FOSTA resulted in a 3.9% increase in sexual assualt rates among sex-workers.

On the other hand, our $R^2$ level is quite low, at $R^2 = 0.012$. This means that only about one percent of the change in outcome — sexual assault rates — can be explained by movement in treatment. However, the $R^2$ value only indicates the goodness of fit, so small $R^2$ values are not necessarily a problem, so long as our $P$ value is below the significance level.

| | Robbery | Assault | Sexaul Assault |
|---|---|---|---|
| | (1) | (2) | (3) |
| treat | −0.010 | 0.079*** | 0.023* |
| | (0.034) | (0.028) | (0.012) |
| post | −0.034 | 0.020 | −0.006 |
| | (0.035) | (0.029) | (0.012) |
| treat:post | 0.045 | −0.037 | 0.039** |
| | (0.049) | (0.041) | (0.017) |
| Constant | 0.328*** | 0.156*** | 0.011 |
| | (0.025) | (0.020) | (0.009) |
| Standard Errors | Standard | Standard | Standard |
| Observations | 1,445 | 1,499 | 1,470 |
| $R^2$ | 0.001 | 0.007 | 0.020 |
| Adjusted $R^2$ | −0.001 | 0.005 | 0.018 |
| Residual Std. Error | 0.465 (df = 1441) | 0.396 (df = 1495) | 0.163 (df = 1466) |
| F Statistic | 0.423 (df = 3; 1441) | 3.377** (df = 3; 1495) | 10.009*** (df = 3; 1466) |
| *Note:* | | | *$p<0.1$; **$p<0.05$; ***$p<0.01$ |

*Table 3: Results of difference-in-differences test across three outcome variables*

To further ascertain the legitimacy of our findings and the parallel trend assumption, I created three separate trend plots with annual means and standard deviations used for the difference-in-differences analysis across both treatment and control groups. From the simple means plots, we see that the parallel trend assumption is plausible for robbery, hard to tell for assault, and rather implausible for sexual assault. This is quite unfortunate, because estimates of causal effects are only as good as the plausibility of the identification assumption, which is the parallel trend assumption.

To more systematically and quantitatively verify the validity of the parallel trend assumption, or lack thereof, I decided to run a placebo test using previous periods. Specifically, I selected 2016 as the dummy treatment year because 2016 Q2 is the point where the treament and control trend lines crossed over in the sexual assualt outcome plot.



*Figure 4: Parallel trends using pre-treatment data, mean and standard deviation trends plot across Robbery, Assault, and Sexual Assualt outcome variables*

Table 4 details the result of running a placebo difference-in-differences test with 2016 as the dummy treatment year. If the parallel-trends assumption were violated in the original evaluated time-range, we would expect to see that the results are statistically significant. This would be very bad for the real difference-in-difference results, because it means that the difference in means would have been statistically significant even without the presence of the intervention — In the context of our analysis, this would mean that sexual assault rates among sex workers would skyrocket in the timeframe even without SESTA-FOSTA. On the flip side, the absence of a significant result in the placebo test means that, at least in our analyzed timeframe, means that the parallel trends assumption could be upheld for the pre-treatment period analyzed.

As shown in table 4, luckily, none of the placebo results across the three outcome variables are significant. This means that although the previously graphed trends plot do not visually look parallel, they are statistically parallel enough to uphold the difference-in-differences result, and the significant findings from table 3.

|  | Robbery Dummy | Assault Dummy | Sexaul Assault Dummy |
|---|---|---|---|
|  | (1) | (2) | (3) |
| treat | −0.042 | 0.002 | 0.006 |
|  | (0.028) | (0.023) | (0.008) |
| post | −0.072** | −0.023 | −0.001 |
|  | (0.032) | (0.025) | (0.009) |
| treat:post | 0.031 | 0.077** | 0.016 |
|  | (0.044) | (0.036) | (0.012) |
| Constant | 0.399*** | 0.179*** | 0.012** |
|  | (0.020) | (0.016) | (0.005) |
| Standard Errors | Standard | Standard | Standard |
| Observations | 1,939 | 1,990 | 2,018 |
| $R^2$ | 0.004 | 0.005 | 0.004 |
| Adjusted $R^2$ | 0.003 | 0.003 | 0.002 |
| Residual Std. Error | 0.478 (df = 1935) | 0.389 (df = 1986) | 0.134 (df = 2014) |
| F Statistic | 2.901** (df = 3; 1935) | 3.115** (df = 3; 1986) | 2.609** (df = 3; 2014) |
| *Note:* |  |  | *p<0.1; **p<0.05; ***p<0.01 |

*Table 4: Placebo difference-in-differences test with 2016 as the dummy treatment year*

*Results across Race*

In qualitative surveys and interviews previously cited in the literature review, one common remark is the disparate impact of violence when categorized by the race of sex workers. Overwhelming, individual testimonies show that white sex-workers face fewer run-ins with the law, and occupy the higher and more protected strata of sex-work. As a result, I evaluated difference-in-differences results of SESTA-FOSTA on sexual assault rates across individual ethnic and racial groups. Although the NCVS survey provides fifteen categories for race, only the first five — white, black, Native American and Alaskan Native, Asian, and Pacific Islander — have any entries across either treatment or control in the evaluated timescope. Out of the five categories, only the first three have any results in the treated group.

Table 4 details the result of difference-in-differences across racial groups.

| | White only | Black only | Native American, Alaskan Native |
|---|---|---|---|
| | (1) | (2) | (3) |
| treat | 0.020 | −0.000 | 0.000 |
| | (0.015) | (0.024) | (0.069) |
| post | 0.006 | 0.002 | 0.000 |
| | (0.015) | (0.025) | (0.072) |
| treat:post | 0.030 | 0.015 | 0.500*** |
| | (0.022) | (0.035) | (0.121) |
| Constant | 0.017 | 0.015 | −0.000 |
| | (0.011) | (0.017) | (0.051) |
| Standard Errors | Standard | Standard | Standard |
| Observations | 1,087 | 255 | 39 |
| $R^2$ | 0.012 | 0.003 | 0.473 |
| Adjusted $R^2$ | 0.010 | −0.009 | 0.428 |
| Residual Std. Error | 0.183 (df = 1083) | 0.140 (df = 251) | 0.169 (df = 35) |
| F Statistic | 4.477*** (df = 3; 1083) | 0.227 (df = 3; 251) | 10.470*** (df = 3; 35) |

Note: $^{*}p<0.1$; $^{**}p<0.05$; $^{***}p<0.01$

*Table 4: difference-in-differences analysis on sexaul assault rates across racial groups*

Notably, while White and Black only categories have no significant results, the Native American and Pacific Islander subgroup is significant at the 0.01 significance level. In fact, the ATT indicates that SESTA-FOSTA's implementation resulted in a fifty-percent increase in sexual rates among Native American, Alaskan Native sex-workers. This is very troubling to the internal validity of our findings —

such a strong disparate impact might indicate that prior results come from the Native American Alaskan Native subgroup even with matching. Moreover, this could potentially indicate that some hidden factor that specifically targets the Native American subgroup is responsible for the change in our outcome variable. On the other hand, the specific finding also doesn't generalize well, with only a sample size of $N = 39$. This means that further robustness tests are needed to determine whether the prior criteria used for identifying sex-workers are too generous, thus yielding false positives.

*Robustness testing*

To further constrain our indicators of who is a sex-worker, a more strict set of criteria is detailed below to verify the robustness of our findings. If significant findings still exist with this new set of indicator variables, then we can be more confident of the internal validity of prior results. Otherwise, our prior statistical significance is probably due to the prevalence of false positives in our dataset. In the chart below, light yellow highlights indicate a tightening of thresholds.

| | | |
|---|---|---|
| **IF (V2014 = 2 or 3)** | Rented unit for cash or no cash rent | **AND** |
| (V2020 = 2, 3, 4, 5, 6, 7, 8, 9, 10) | Type of living quarters is non-transient hotel, motel, transient hotel, motel, rooming house, mobile home or trailer, boarding house | **AND** |
| (V2024 != 1, 2, 98, 99) | Number of units in structure is greater than 2, and not residue or out of universe | **AND** |
| (V2025A != 2) | Not living in a gated/walled community | **AND** |
| (V2025B != 2) | Building does not have restricted access | **AND** |
| (16 <= V2033 <= 50) | Aged between 16 to 50 | **AND** |
| (V2036 = 2 or trans = TRUE) | Either female or trans-identifying | **AND** |
| (V2122 contains "lone") | Living alone partner in family structure code | **AND** |
| (V3033 > 2) | Moved more than twice in the past five years | **AND** |
| (V2026 <50000) | Annual household income less than $50,000 | **AND** |
| (V2074 = 1 and V2075 = 2) | Operate business from address but does not have a sign | **THEN** |
| Treat = 1 | Respondent is categorized as a sex-worker | |

One caveat encountered when trying to replicate the procedure for the robustness test is that the subsetted dataset after matching returned no positive outcome units in the analyzed time frame (from 2017Q1 to 2017Q4, and then from 2019Q1 to 2019Q4). This means that we could not discern a potential lack of statistical significance as due to the actual treatment, or due to the matching procedure. As a result, I did not match on any confounders, and opted instead to add the three covariates in the regression equation used for difference-in-differences.

The result of the robustness test is shown in table 5 below. As we can see from the results, sadly, the significant findings in the prior analysis did not survive our robustness test. We now see virtually zero percent of movement in the outcome can be explained by the treatment. This means that our earlier findings are likely due to variance in the dataset, rather than approaching population-level ground truth.

| | Sexaul Assault |
|---|---|
| treat | −0.011 |
| | (0.057) |
| post | 0.001 |
| | (0.002) |
| education | −0.0003 |
| | (0.0002) |
| urban.rural | −0.007*** |
| | (0.002) |
| race | −0.0004 |
| | (0.001) |
| treatTRUE:post | −0.00000 |
| | (0.093) |
| Constant | 0.031*** |
| | (0.006) |
| Standard Errors | Standard |
| Observations | 21,080 |
| $R^2$ | 0.0005 |
| Adjusted $R^2$ | 0.0002 |
| Residual Std. Error | 0.128 (df = 21073) |
| F Statistic | 1.661 (df = 6; 21073) |
| Note: | *p<0.1; **p<0.05; ***p<0.01 |

*Table 5, difference in differences without matching on Sexual Assault, robustness check*

**Discussion and Critique**

There are numerous limitations and potential areas for improvement in this research design. The most significant limitation is the difficulty and uncertainty in determining who is a sex-worker in the NCVS survey results. Sex-work as a profession has long been considered taboo, and there exists significant overlap between sex-workers and victims of human trafficking. As a result, respondents of even targeted surveys often do not feel comfortable or safe to self-identify as a sex-worker. With NCVS's dataset, the greatest challenge is to encapsulate sex-workers' diverse backgrounds within the limited datafields and response bias. For instance, exclusive escorts may share an incredibly different economic and residential profile as victims of human-trafficking, who often live with other victims and their trafficker and lack any disposable income.

Notably, while qualitative interviews cited in this paper tend to overestimate the impact of SESTA-FOSTA on sex-workers, quantitative, wide-casting survey results in the NCVS tend to underestimate the impact. This is because smaller-scale, community-based organizations obtain their respondents from those obtaining their services, who often happen to be sex-workers in the periphery with dire needs such as housing insecurity and drug abuse. On the other hand, NCVS-style surveys find their respondents based on landlines and permanent addresses, which many of those most vulnerable groups lack. As a result, each method of data collection is biased and flawed in their own way.

A potential way to improve the design is with better datasets and data collection in the future. For instance, the NCVS questionnaire could incorporate some of the screening questions from New York Trafficking Assessment Project's toolkit. Questions such as whether the respondent personally has access to their legal documents, whether they are allowed to leave or take breaks from their work, and whether non-consensual photos were taken of respondents during work retain respondent anonymity but effectively reveal the occupation of respondents. Moreover, partnerships with community-based organizations could also yield a more complete census demographics.
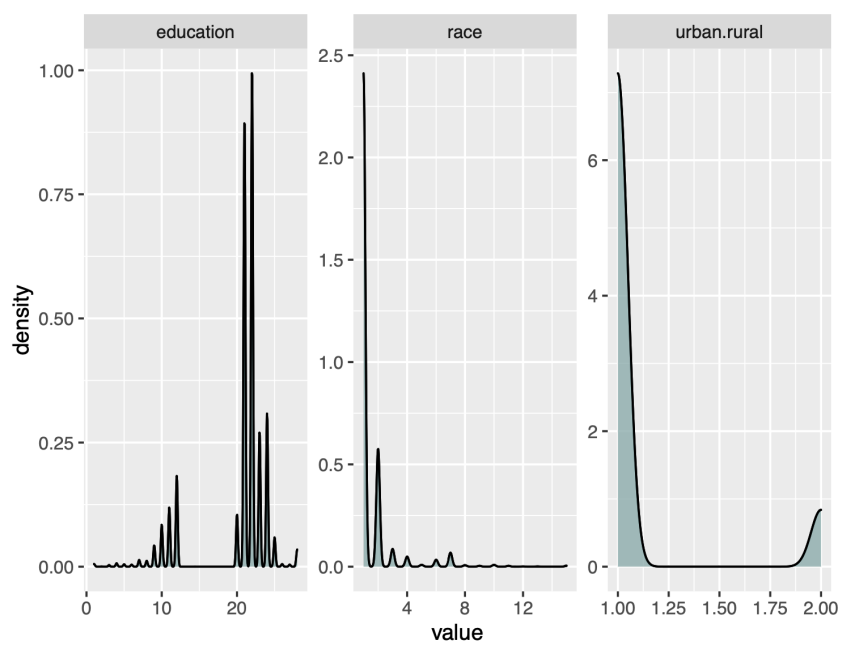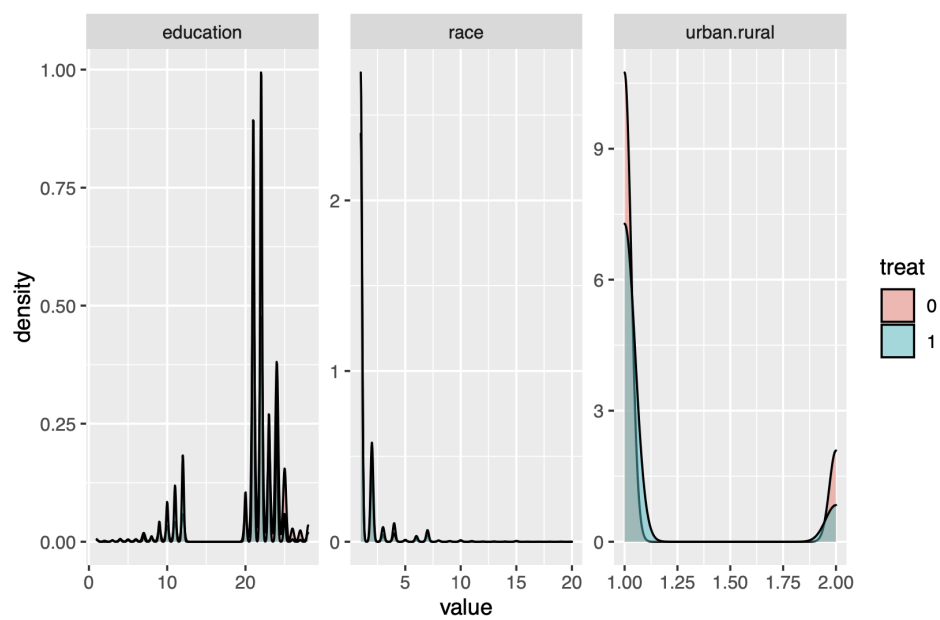
Besides the limitations in identifying the treatment group, a common issue in studying crimes against women and sex-workers is the underreporting of rape and sexual assault. According to the Justice Department, nearly eighty percent of rapes and sexual assualts go unreported. Although this phenomena should be equally prominent in both the treated and control groups, it marks the limitations of our analysis: Even the most rigorous causal analysis only reveals the difference between the treated and control groups, rather than the true counts of victimization.

For the specific methods used in this paper, the biggest threats to internal validity are the issue of false positives and self-selection of the treatment group before and after 2018. As revealed in the robustness test, identifying sex-workers based on demographics traits present in the NCVS results lead to significant overcounting, which resulted in us first finding a statistically significant result when there probably isn't one in the given dataset.

Additionally, respondents may change status from sex-worker to non-sex-worker, and vice versa, during the years evaluated in the study. For instance, the 2019 coronavirus pandemic resulted in a massive loss of employment for most in-person sex workers. From research conducted by the Social Sciences Research Council, sex workers often lost their primary source of income because of the pandemic, and must seek work and financial stability through other means. This means that the treatment and control groups may have overlapping responses, as people switch occupations in a predictable way. Luckily, for the purpose of this research, we can track each respondent's NCVS ID, and 2017 to 2019 fall within one full rotation cycle. By tracing for whether the same respondent appears in the treated and control groups before or after the SESTA-FOSTA enactment, we can remove the respondents who switched their treatment status.

Lastly, this research project does not differentiate between consentual sex-workers and victims of sexual trafficking, although a fundemental difference exist betweent the two overlapping groups. Though lawmakers often draft legislations to target human trafficking, all sex-workers — regardless of whther they are trafficked — are impacted by legislations such as SESTA-FOSTA.

# Appendix

**Endnotes**

1. National Institute of Justice, Michael Shively, Kristina Kliorys, Kristin Wheeler, and Dana Hunt, A National Overview of Prostitution and Sex Trafficking Demand Reduction Efforts, Final Report §. 238796 (2012). https://www.ojp.gov/pdffiles1/nij/grants/238796.pdf. Award Number: 2008-IJ-CX-0010

2. National Institute of Justice, Shively et. all

3. D Blunt and A Wolf, 'Erased: The impact of FOSTA-SESTA and the removal of Backpage on sex workers', Anti-Trafficking Review, issue 14, 2020, pp. 117-121, https://doi.org/10.14197/atr.201220148

4. Mullin, Joe. "How Fosta Could Give Hollywood the Filters It's Long Wanted." Electronic Frontier Foundation. Electronic Frontier Foundation, March 16, 2018. https://www.eff.org/deeplinks/2018/03/how-fosta-will-get-hollywood-filters-theyve-long-wanted.

5. "Screening for Human Trafficking Guidelines for Administering the Trafficking Victim Identification Tool (TVIT)." New York: Vera Institute of Justice, June 2014. NCJ #246713. https://www.vera.org/downloads/publications/human-trafficking-identification-tool-and-user-guidelines.pdf

6. United Nations Office on Drugs And Crime, HUMAN TRAFFICKING INDICATORS §. Accessed March 27, 2022. https://www.unodc.org/pdf/HT_indicators_E_LOWRES.pdf.

7. D Blunt and A Wolf

8. Government Accountability Office, SEX TRAFFICKING: Online Platforms and Federal Prosecutions § (2021). https://www.gao.gov/assets/gao-21-385.pdf.

9. Romano, Aja. "A New Law Intended to Curb Sex Trafficking Threatens the Future of the Internet as We Know It." Vox. Vox, April 13, 2018. https://www.vox.com/culture/2018/4/13/17172762/fosta-sesta-backpage-230-internet-freedom.

10. D Blunt and A Wolf

11. Gallay, Amelia. "Sex Sells, but Not Online: Tracing the Consequences of Fosta-Sesta." BJCL. Berkeley Journal of Criminal Law, December 4, 2021. https://www.bjcl.org/blog/sex-sells-but-not-online-tracing-the-consequences-of-fosta-sesta.

12. Whitford, Emma. "There's No Such Thing as a Low-Level Arrest When You're Undocumented." Jezebel. G/O Media Inc., December 19, 2018. https://jezebel.com/theres-no-such-thing-as-a-low-level-arrest-when-youre-u-1831205673.

13. Bureau of Justice Statistics, National Crime Victimization Survey (NCVS) § (2020). https://bjs.ojp.gov/data-collection/ncvs#surveys-0.

14. National Commission on COVID-19 and Criminal Justice, Richard Rosenfeld, and Ernesto Lopez , Pandemic, Social Unrest, and Crime in U.S. Cities: March 2021 Update § (2021). https://counciloncj.org/wp-content/uploads/2021/11/Pandemic_Social_Unrest_and_Crime_in_US_Cities_-_March_2021_Update.pdf.

15. National Commission on COVID-19 and Criminal Justice, Richard Rosenfeld, and Ernesto Lopez

Yichuan Shi Project Data-cleaning Script

## Cleanup

The purpose of the cleanup script is to independently clean and merge the three test data files into a single file for running the difference-in- differences analysis

First, load and isolate data with fields of interest

```
library(stringi)

load("test3.rda")

# columns of interest:
c.interest <- c(
  "YEARQ",
  "YEAR", # year and quarter of interview
  "V2006", # household number
  "V2008", # panel and rotation group
  "V2014", # tenure, owned house/cash rent/no cash rent
  "V2016", # land use, urban/rural/residue
  "V2020", # type of living quarters
  "V2024", # number of housing units in structure
  "V2025", # direct outside access
  "V2025A", # gated or walled community
  "V2025B", # building with restricted access
  "V2026", # household income
  "V2032", # principle person relation to ref person
  "V2033", # principle person age
  "V2036", # principle person sex
  "V2038", # principle person educational attainment
  "V2040A", # principle person race recode
  "V2073", # no. crime incident reports
  "V2074", # operate business from address
  "V2075", # sign on premises for business
  "V2105A", # targeted because of race
  "V2105C", # targeted because of ethnicity
  "V2105E", # targeted because of gender
  "V2105F", # targeted because of sexuality
  "V2122", # family structure code
  "V3033", # times moved in the last 5 years
  "V3046", # forced or coerced unwanted sex
  "V3047", # no. times unwanted sex
  "V3048", # call police to report something else
  "V3053", # no. times attack, threat, theft
  "V3054", # crime not reported to police
  # started since 2017: -------------
  "V3083", # citizenship status
  "V3085", # gender identity at birth
  "V3086", # current gender identity
  # end of 2017 only block --------------
  "V3074", # which best describes your job
  "V3075", # is employment private, govt, or self
  "V4093", # how attacked 1=at least 1 good entry in codes 1-14
  "V4482B", # collapsed occupation code
```

```r
  "V4484", # incident occur at work site
  "V4529" # Type of crime code new
)

intersect <- intersect(names(da38136.0003), c.interest)
s <- da38136.0003[,intersect]
# s <- s[s$YEARQ>2007.4,]

s <- s[!is.na(s$YEARQ),] # get rid of na entries

# change year field into quarterly
# ie 2007.1 = 2007
   #2007.2 = 2007.25
   #2007.3 = 2007.50
   #2007.4 = 2007.75

s$YEARQ <- as.character(s$YEARQ)

# first quarter
s$YEARQ <- gsub("\\b.1\\b", "", s$YEARQ)

# second quarter
s$YEARQ <- gsub("\\b.2\\b", ".25", s$YEARQ)

# third quarter
s$YEARQ <- gsub("\\b.3\\b", ".50", s$YEARQ)

# fourth quarter
s$YEARQ <- gsub("\\b.4\\b", ".75", s$YEARQ)

# back to numeric:
s$YEARQ <- as.numeric(s$YEARQ)
```

Now, we want to make our outcome variables: robbery, personal assault, and sexual assault

```r
# robery

rob <- c(
  "\\<(05)\\>",
  "\\<(06)\\>",
  "\\<(07)\\>",
  "\\<(08)\\>",
  "\\<(09)\\>",
  "\\<(10)\\>",
  "\\<(31)\\>",
  "\\<(32)\\>",
  "\\<(33)\\>"
)

s$robbery <- grepl(paste(rob, collapse = "|"), s$V4529)
# convert from boolean to integer
s$robbery <- as.integer(s$robbery)

# assault, not including assault with sexual attack
```

```r
assault <- c(
  "\\<(11)\\>",
  "\\<(12)\\>",
  "\\<(13)\\>",
  "\\<(14)\\>",
  "\\<(17)\\>",
  "\\<(20)\\>"
)

s$assault <- grepl(paste(assault, collapse = "|"), s$V4529)
# convert from boolean to integer
s$assault <- as.integer(s$assault)

# sexual assault, including verbal attack and serious physical injury
sexual.assault <- c(
  "\\<(01)\\>",
  "\\<(02)\\>",
  "\\<(03)\\>",
  "\\<(04)\\>",
  "\\<(15)\\>",
  "\\<(16)\\>",
  "\\<(18)\\>",
  "\\<(19)\\>"
)

s$sexual.assault <- grepl(paste(sexual.assault, collapse = "|"), s$V4529)
# convert from boolean to integer
s$sexual.assault <- as.integer(s$sexual.assault)
```

Now, create our treatment variable using compilation of indicators for potential sex workers

```r
# living quarters
living.quarters <- c(
  "\\<(01)\\>",
  "\\<(02)\\>",
  "\\<(03)\\>",
  "\\<(04)\\>",
  "\\<(05)\\>",
  "\\<(06)\\>",
  "\\<(07)\\>",
  "\\<(08)\\>",
  "\\<(09)\\>",
  "\\<(10)\\>"
)

number.units <- c(
  "\\<(98)\\>",
  "\\<(99)\\>"
)

income <- c(
  "\\<(13)\\>",
  "\\<(14)\\>",
  "\\<(15)\\>",
```

```r
  "\\<(16)\\>",
  "\\<(17)\\>",
  "\\<(18)\\>",
  "\\<(99)\\>"
)

# trans indicator
s$V3085 <- as.character(s$V3085)
s$V3086 <- as.character(s$V3086)

trans.discount <- c(
  "\\b(3)\\b",
  "\\b(8)\\b",
  "\\b(-1)\\b",
  "\\b(9)\\b"
)

s$trans <-
  (!is.na(s$V3085) & !is.na(s$V3086)) &
  (s$V3085 != s$V3086) &
  !grepl(paste(trans.discount,
               collapse = "|"), s$V3085) &
  !grepl(paste(trans.discount,
               collapse = "|"), s$V3086)

# map the treatment indicator

s$treat <-
  ((s$V2014=="(2) Rented for cash" | s$V2014=="(3) No cash rent") &
  grepl(paste(living.quarters, collapse = "|"), s$V2020) & # 49771
  !grepl(paste(number.units, collapse = "|"), s$V2024) & # 49771
  s$V2025A == "(2) No" & # not gated/walled community 42813
  s$V2025B == "(2) No" & # no restricted access 39760
  (s$V2033 > 15 & s$V2033 < 61) & # age range 36436
  (as.character(s$V2036)=="(2) Female"| s$trans == TRUE) & # either female or transgender 25923
  grepl("Lone", s$V2122) &  # lone male/female
  s$V3033>2 & # moved more than twice in past 5 years 13613
  !grepl(paste(income, collapse = "|"), s$V2026)) |  # income less than 50000
  (s$V2074 == "(01) Yes" & s$V2075 == "(02) No" ) # or operate business but no sign)

# to numeric
s$treat <- as.integer(s$treat)

# remove na entries
s <- s[!is.na(s$treat),]
```

Now let's run some sanity checks to see if the numbers make sense at all

```r
library(ggplot2)
library(gridExtra)

# simple average per year of treat and control groups across outcome variables
t <- s[s$treat == 1,]
c <- s[s$treat == 0,]
```

```r
mean.assault.t <- aggregate(t$assault, by=list(Year = t$YEARQ), FUN=mean)
mean.assault.c <- aggregate(c$assault, by=list(Year = c$YEARQ), FUN=mean)

mean.rob.t <- aggregate(t$robbery, by=list(Year = t$YEARQ), FUN=mean)
mean.rob.c <- aggregate(c$robbery, by=list(Year = c$YEARQ), FUN=mean)

mean.sa.t <- aggregate(t$sexual.assault, by=list(Year = t$YEARQ), FUN=mean)
mean.sa.c <- aggregate(c$sexual.assault, by=list(Year = c$YEARQ), FUN=mean)

mean.assault.c <- subset(mean.assault.c, Year > 2003.75)
mean.rob.c <- subset(mean.rob.c, Year > 2003.75)
mean.sa.c <- subset(mean.sa.c, Year > 2003.75)

mean.assault.c <- subset(mean.assault.c, Year < 2018)
mean.rob.c <- subset(mean.rob.c, Year < 2018)
mean.sa.c <- subset(mean.sa.c, Year < 2018)

mean.assault.t <- subset(mean.assault.t, Year < 2018)
mean.rob.t <- subset(mean.rob.t, Year < 2018)
mean.sa.t <- subset(mean.sa.t, Year < 2018)

# plot
p.assault <- ggplot(NULL, aes(x=Year, y=x)) +
        geom_point(data = mean.assault.c, color="blue") +
        geom_point(data = mean.assault.t, color="red") +
        geom_smooth(data = mean.assault.c, method="lm", colour = "blue") +
        geom_smooth(data = mean.assault.t, method="lm", colour = "red") +
        ggtitle("Assault pre 2018") +
        geom_vline(xintercept = 2018) +
        scale_color_manual(name='Legend',
                breaks=c('Control', 'Treatment'),
                values=c('Control'='blue', 'Treatment'='red'))

p.robery <- ggplot(NULL, aes(x=Year, y=x)) +
        geom_point(data = mean.rob.c, color="blue") +
        geom_point(data = mean.rob.t, color="red") +
        geom_smooth(data = mean.rob.c, method="lm", colour = "blue") +
        geom_smooth(data = mean.rob.t, method="lm",colour = "red") +
        ggtitle("Robbery pre 2018") +
        geom_vline(xintercept = 2018)


p.sexual.assault <- ggplot(NULL, aes(x=Year, y=x)) +
        geom_point(data = mean.sa.c, color="blue") +
        geom_point(data = mean.sa.t, color="red") +
        geom_smooth(data = mean.sa.c, method="lm", colour = "blue") +
        geom_smooth(data = mean.sa.t, method="lm", colour = "red") +
        ggtitle("Sexual Assault pre 2018") +
        geom_vline(xintercept = 2018)

grid.arrange(p.assault, p.robery, p.sexual.assault, ncol=2)

## `geom_smooth()` using formula 'y ~ x'
```
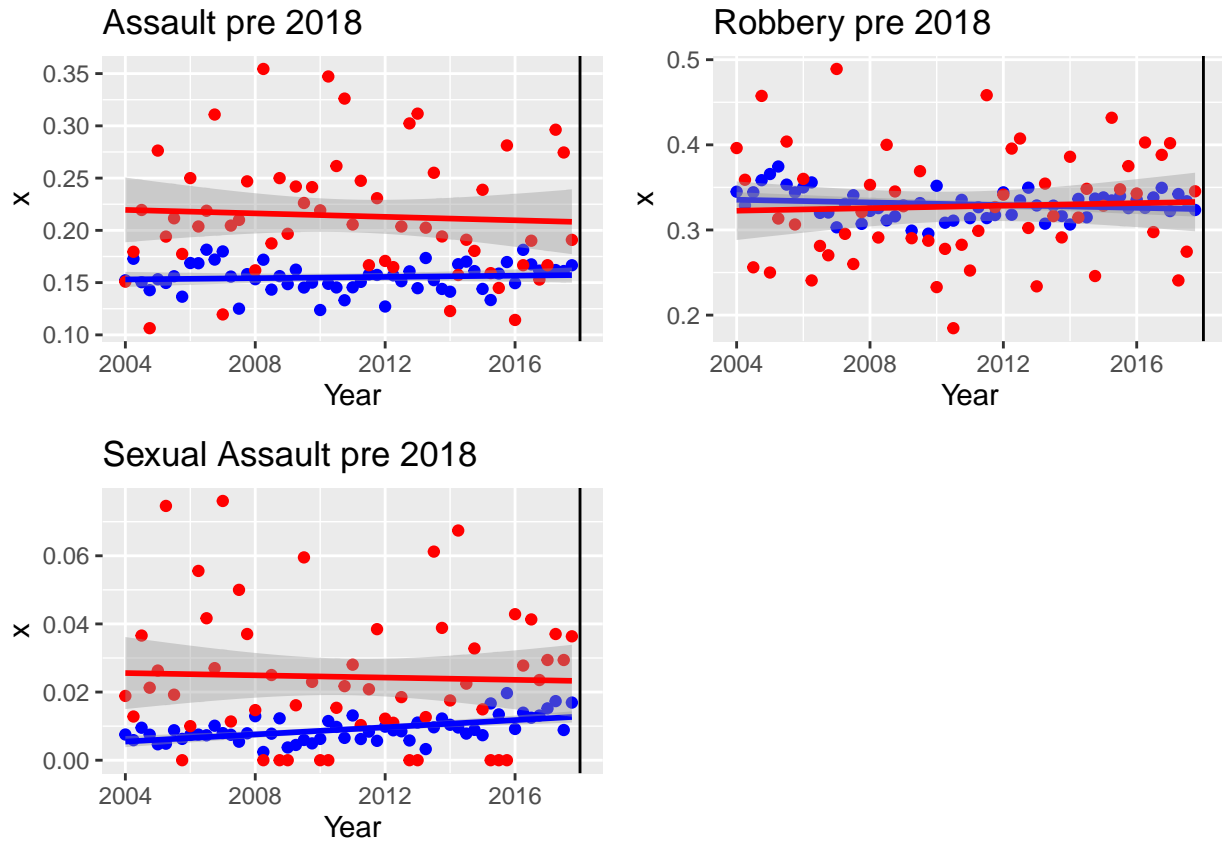
```
## `geom_smooth()` using formula 'y ~ x'
## `geom_smooth()` using formula 'y ~ x'
## `geom_smooth()` using formula 'y ~ x'
## `geom_smooth()` using formula 'y ~ x'
## `geom_smooth()` using formula 'y ~ x'
```

### Assault pre 2018



### Robbery pre 2018



### Sexual Assault pre 2018



Obviously, there is not a lot of parallel trends in some of our outcomes of interest. However, we have not yet controlled for potential covariates in our study. Before doing the diff-in-diff exercise, we should identify the covariates to match for: Education, race, ethnicity, gender, sexuality

```
# save cleaned data
s <- s[, c(
  "YEARQ",
  "YEAR",
  "treat",
  "robbery",
  "assault",
  "sexual.assault",
  "V2038",
  "V2016",
  "V2040A"

)]
save(s, file="cleaned.RData")
```

# remember to cut control data and run placebo test

This file pre-processes our data set via matching with the identified covariates:

"V2038", principle person educational attainment
"V2016", urban/rural/residue
"V2040A", principle person race

```
load("cleaned.RData")

library(Matching)
library(ebal)
library(ggplot2)
library(dplyr)
library(xtable)
library(tidyr)
library(foreign)

s.no.na <- s[!is.na(s$V2040A),]
s.no.na <- s.no.na[!is.na(s.no.na$V2016),]
s.no.na <- s.no.na[!is.na(s.no.na$V2038),]

# convert categorical to numerical
s.no.na$V2038 <- unclass(s.no.na$V2038)
s.no.na$V2016 <- unclass(s.no.na$V2016)
s.no.na$V2040A <- unclass(s.no.na$V2040A)

colnames(s.no.na) <- c(
  "year.q",
  "year",
  "treat",
  "robbery",
  "assault",
  "sexual.assault",
  "education",
  "urban.rural",
  "race"
)
```

To actually check the match balance
```
# robbery
SW.formula.rb <- formula(treat ~ education + urban.rural + race)

mb.unmatched.rb <- MatchBalance(SW.formula.rb, data = s.no.na, print.level = 0)

varnames.rb <- c("education", "urban.rural", "race")

t1 <- baltest.collect(mb.unmatched.rb,
                var.names = varnames.rb,
                after = F) %>%
  .[,c(1,2,6)]
  #xtable(caption = "Balance Pre-Matching robbery") %>%
  #print(., include.rownames=T, caption.placement="top", comment = F)

# assault
SW.formula.as <- formula(treat ~ education + urban.rural + race)
```

```r
mb.unmatched.as <- MatchBalance(SW.formula.as, data = s.no.na, print.level = 0)
varnames.as <- c("education", "urban.rural", "race")
t2 <- baltest.collect(mb.unmatched.as, var.names = varnames.as,after = F) %>%
  .[,c(1,2,6)]
  #xtable(caption = "Balance Pre-Matching Assault") %>%
  # print(., include.rownames=T, caption.placement="top", comment = F)

# sexual assault
SW.formula.sa <- formula(treat ~ education + urban.rural + race)
mb.unmatched.sa <- MatchBalance(SW.formula.sa, data = s.no.na, print.level = 0)
varnames.sa <- c("education", "urban.rural", "race")
t3 <- baltest.collect(mb.unmatched.sa, var.names = varnames.sa, after = F) %>%
  .[,c(1,2,6)]
  #xtable(caption = "Balance Pre-Matching Sexual Assault") %>%
  #print(., include.rownames=T, caption.placement="top", comment = F)

# make into dataframe
t1 <- as.data.frame(t1)
t2 <- as.data.frame(t2)
t3 <- as.data.frame(t3)

tpre <- rbind(t1, t2, t3)
rownames(tpre) <- c(
  "education.robbery",
  "urban.rural.robbery",
  "race.robbery",
  "education.assault",
  "urban.rural.assault",
  "race.assault",
  "education.sexual.assault",
  "urban.rural.sexual.assault",
  "race.sexual.assault"

)

# make into talbe
xtable(tpre, caption = "Balance Pre-Matching") %>%
  print(., include.rownames=T, caption.placement="top", comment = F)
```

Table 1: Balance Pre-Matching

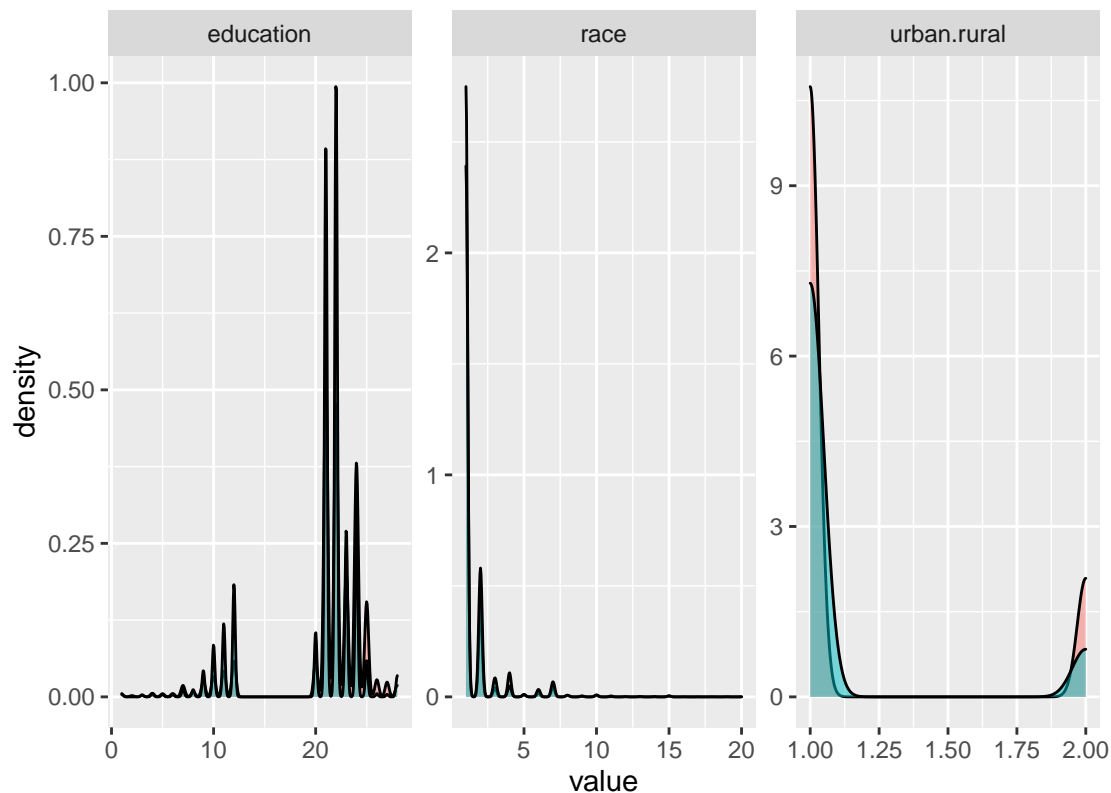|  | mean.Tr | mean.Co | T pval |
|---|---|---|---|
| education.robbery | 20.33 | 21.29 | 0.00 |
| urban.rural.robbery | 1.10 | 1.16 | 0.00 |
| race.robbery | 1.55 | 1.44 | 0.00 |
| education.assault | 20.33 | 21.29 | 0.00 |
| urban.rural.assault | 1.10 | 1.16 | 0.00 |
| race.assault | 1.55 | 1.44 | 0.00 |
| education.sexual.assault | 20.33 | 21.29 | 0.00 |
| urban.rural.sexual.assault | 1.10 | 1.16 | 0.00 |
| race.sexual.assault | 1.55 | 1.44 | 0.00 |

This is very bad– all of our covariates seem statistically significant. We need to match :(

But first, let's look at density comparison

```
s.no.na %>%
  dplyr::select(education,
                urban.rural,
                race,
                treat) %>%
  mutate(treat = as.factor(treat)) %>%
  gather(key= variable, value = value, -treat) %>%
  ggplot(aes(value, fill=treat))+
  geom_density(alpha=0.5)+
  facet_wrap(~variable, scales="free")+
  ggtitle("Pre-treatment balance of covariates")
```

```
## Warning: attributes are not identical across measure variables;
## they will be dropped
```



Pre−treatment balance of covariates

Let's now actually implement matching and check our covariate balance again

```
# robbery
match.rb <- Match(Y = s.no.na$robbery,
                  Tr = s.no.na$treat,
                  X = s.no.na[,varnames.rb],
                  Weight = 2,
                  ties=FALSE)
matched.rb <- MatchBalance(SW.formula.rb,
                           match.out = match.rb,
                           data = s.no.na,
                           print.level = 0)
```

```r
t1 <- baltest.collect(matched.rb,
                  var.names = varnames.rb, after=TRUE) %>%
  .[,c(1,2,6)]
  #xtable(caption="Post-matching balance") %>%
  #print(., include.rownames=T, caption.placement="top", comment = F)

# assault
match.as <- Match(Y = s.no.na$assault,
                  Tr = s.no.na$treat,
                  X = s.no.na[,varnames.as],
                  Weight = 2,
                  ties=FALSE)
matched.as <- MatchBalance(SW.formula.as,
                           match.out = match.as,
                           data = s.no.na,
                           print.level = 0)

t2 <- baltest.collect(matched.as,
                  var.names = varnames.as, after=TRUE) %>%
  .[,c(1,2,6)]
  #xtable(caption="Post-matching balance Assault") %>%
  #print(., include.rownames=T, caption.placement="top", comment = F)

# sexual assault
match.sa <- Match(Y = s.no.na$sexual.assault,
                  Tr = s.no.na$treat,
                  X = s.no.na[,varnames.sa],
                  Weight = 2,
                  ties=FALSE)
matched.sa <- MatchBalance(SW.formula.sa,
                           match.out = match.sa,
                           data = s.no.na,
                           print.level = 0)

t3 <- baltest.collect(matched.sa,
                  var.names = varnames.sa, after=TRUE) %>%
  .[,c(1,2,6)]
  #xtable(caption="Post-matching balance Sexual Assault") %>%
  #print(., include.rownames=T, caption.placement="top", comment = F)

# make into dataframe
t1 <- as.data.frame(t1)
t2 <- as.data.frame(t2)
t3 <- as.data.frame(t3)

tpre <- rbind(t1, t2, t3)
rownames(tpre) <- c(
  "education.robbery",
  "urban.rural.robbery",
  "race.robbery",
  "education.assault",
  "urban.rural.assault",
  "race.assault",
```

```
    "education.sexual.assault",
    "urban.rural.sexual.assault",
    "race.sexual.assault"

)

# make into talbe
xtable(tpre, caption = "Balance Post-Matching") %>%
  print(., include.rownames=T, caption.placement="top", comment = F)
```

Table 2: Balance Post-Matching

|  | mean.Tr | mean.Co | T pval |
|---|---|---|---|
| education.robbery | 20.33 | 20.33 | 0.08 |
| urban.rural.robbery | 1.10 | 1.10 | 1.00 |
| race.robbery | 1.55 | 1.55 | 0.08 |
| education.assault | 20.33 | 20.33 | 0.56 |
| urban.rural.assault | 1.10 | 1.10 | 1.00 |
| race.assault | 1.55 | 1.55 | 0.08 |
| education.sexual.assault | 20.33 | 20.33 | 0.08 |
| urban.rural.sexual.assault | 1.10 | 1.10 | 1.00 |
| race.sexual.assault | 1.55 | 1.55 | 0.08 |

Visualize balance across our three outcome variables

```
# robbery
d_matched.rb <- s.no.na[c(match.rb$index.treated, match.rb$index.control),]

d_matched.rb %>%
  dplyr::select(education,
                urban.rural,
                race,
                treat) %>%
  mutate(treat = as.factor(treat)) %>%
  gather(key= variable, value = value, -treat) %>%
  ggplot(aes(value, fill=treat))+
  geom_density(alpha=0.5)+
  facet_wrap(~variable, scales="free")
```
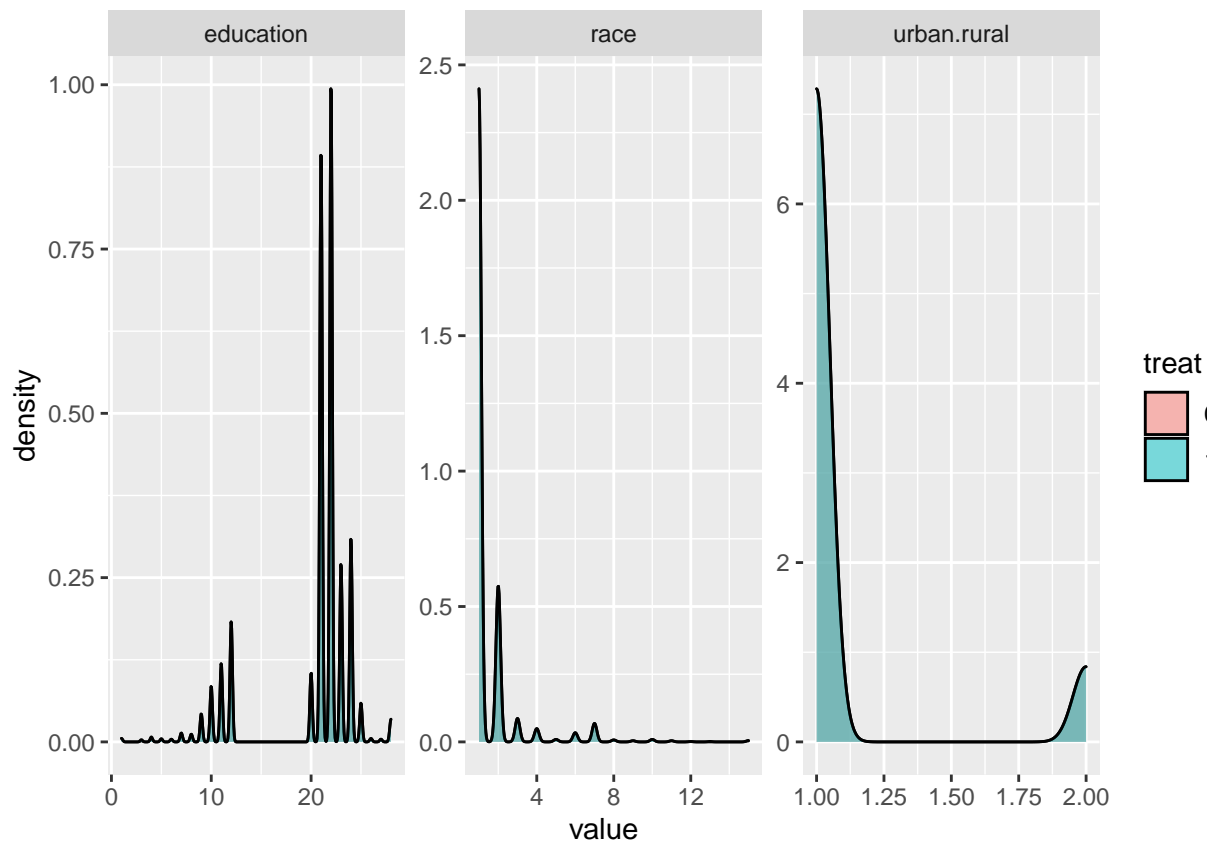
```
# assault
d_matched.as <- s.no.na[c(match.as$index.treated, match.as$index.control),]

d_matched.as %>%
  dplyr::select(education,
                urban.rural,
                race,
                treat) %>%
  mutate(treat = as.factor(treat)) %>%
  gather(key= variable, value = value, -treat) %>%
  ggplot(aes(value, fill=treat))+
  geom_density(alpha=0.5)+
  facet_wrap(~variable, scales="free")
```
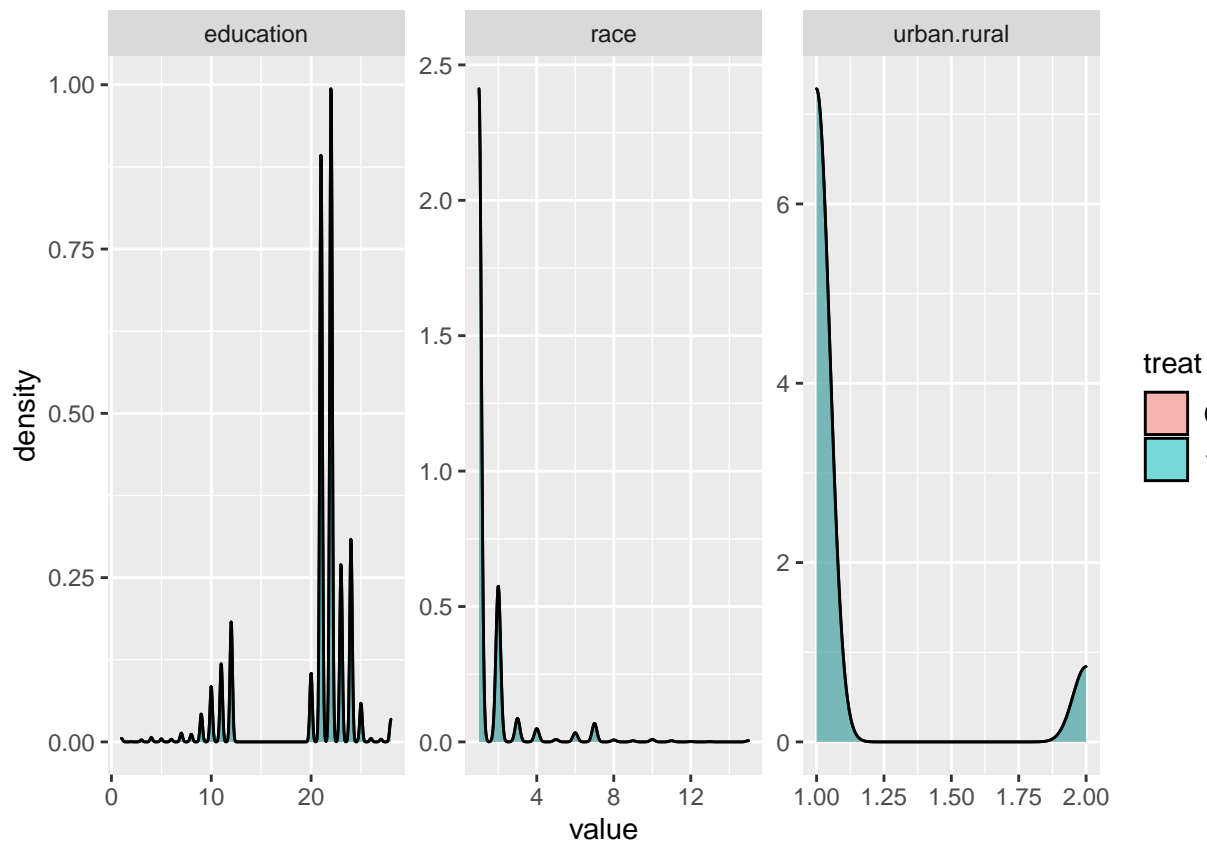
```r
# sexual assault
d_matched.sa <- s.no.na[c(match.sa$index.treated, match.sa$index.control),]
d_matched.sa %>%
  dplyr::select(education,
                urban.rural,
                race,
                treat) %>%
  mutate(treat = as.factor(treat)) %>%
  gather(key= variable, value = value, -treat) %>%
  ggplot(aes(value, fill=treat))+
  geom_density(alpha=0.5)+
  facet_wrap(~variable, scales="free")
```
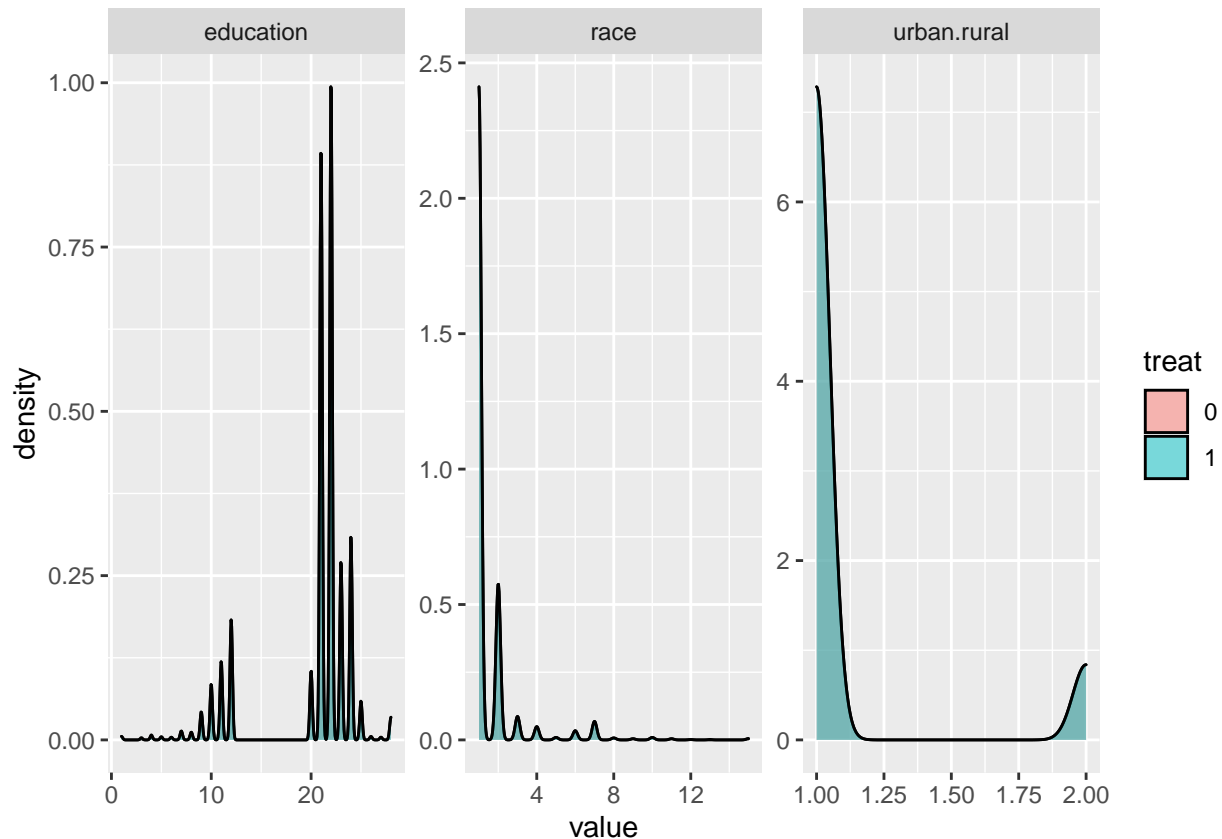
save matched datasets for difference in differences

```
save(d_matched.rb, file="robbery.RData")
save(d_matched.as, file="assault.RData")
save(d_matched.sa, file="sexualAssault.RData")
```

Sexual assualt without matching on race:

```
# sexual assault witout race matching
varnames.nr <- c("education", "urban.rural")
SW.formula.sa.nr <- formula(treat ~ education + urban.rural)
match.sa.nr <- Match(Y = s.no.na$sexual.assault,
                     Tr = s.no.na$treat,
                     X = s.no.na[,varnames.sa],
                     Weight = 2,
                     ties=FALSE)
matched.sa.nr <- MatchBalance(SW.formula.sa.nr,
                           match.out = match.sa.nr,
                           data = s.no.na,
                           print.level = 0)

baltest.collect(matched.sa.nr,
                var.names = varnames.nr, after=TRUE) %>%
  .[,c(1,2,6)] %>%
  xtable(caption="Post-matching balance Sexual Assault, without race") %>%
  print(., include.rownames=T, caption.placement="top", comment = F)

## \begin{table}[ht]
## \centering
```

```
## \caption{Post-matching balance Sexual Assault, without race}
## \begin{tabular}{rrrr}
##   \hline
##  & mean.Tr & mean.Co & T pval \\
##   \hline
## education & 20.33 & 20.33 & 0.56 \\
##   urban.rural & 1.10 & 1.10 & 1.00 \\
##    \hline
## \end{tabular}
## \end{table}
```

```r
d_matched.sa.nr <- s.no.na[c(match.sa.nr$index.treated, match.sa.nr$index.control),]
save(d_matched.sa.nr, file="sa_no_race.RData")
```

Yichuan Shi

Final Project Diff-in-Diff analysis

```r
# load dataset
load("robbery.RData")
load("assault.RData")
load("sexualAssault.RData")

# subset data
r <- subset(d_matched.rb, (year > 2016) & (year < 2020))
a <- subset(d_matched.as, (year > 2016) & (year < 2020))
s <- subset(d_matched.sa, (year > 2016) & (year < 2020))
# remove treatment year
r <- subset(r, year != 2018)
a <- subset(a, year != 2018)
s <- subset(s, year != 2018)

# use regression to estimate difference in difference

r$post <- ifelse(r$year > 2018, 1, 0)
a$post <- ifelse(a$year > 2018, 1, 0)
s$post <- ifelse(s$year > 2018, 1, 0)

# call regression
fit.reg.r <- lm(robbery ~ treat * post, data = r)
summary(fit.reg.r)
```

```
##
## Call:
## lm(formula = robbery ~ treat * post, data = r)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -0.3280 -0.3277 -0.3173  0.6720  0.7064
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.32773    0.02463  13.304   <2e-16 ***
## treat       -0.01042    0.03358  -0.310    0.756
## post        -0.03410    0.03474  -0.982    0.326
## treat:post   0.04477    0.04924   0.909    0.363
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4655 on 1441 degrees of freedom
## Multiple R-squared:  0.0008801,  Adjusted R-squared:  -0.0012
## F-statistic: 0.4231 on 3 and 1441 DF,  p-value: 0.7365
```

```r
fit.reg.a <- lm(assault ~ treat * post, data = a)
summary(fit.reg.a)
```

```
##
## Call:
## lm(formula = assault ~ treat * post, data = a)
##
```

1

```
## Residuals:
##     Min      1Q Median      3Q     Max
## -0.2356 -0.2187 -0.1760 -0.1562  0.8438
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.15617    0.01989   7.853 7.72e-15 ***
## treat        0.07941    0.02780   2.856  0.00435 **
## post         0.01983    0.02853   0.695  0.48723
## treat:post  -0.03676    0.04119  -0.892  0.37234
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.3963 on 1495 degrees of freedom
## Multiple R-squared:  0.006731,   Adjusted R-squared:  0.004738
## F-statistic: 3.377 on 3 and 1495 DF,  p-value: 0.01773
```

```r
fit.reg.s <- lm(sexual.assault ~ treat * post, data = s)
summary(fit.reg.s)
```

```
##
## Call:
## lm(formula = sexual.assault ~ treat * post, data = s)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.06752 -0.03365 -0.01093 -0.00531  0.99469
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.010929   0.008532   1.281   0.2004
## treat        0.022725   0.011698   1.943   0.0522 .
## post        -0.005624   0.011977  -0.470   0.6388
## treat:post   0.039494   0.017122   2.307   0.0212 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1632 on 1466 degrees of freedom
## Multiple R-squared:  0.02007,    Adjusted R-squared:  0.01807
## F-statistic: 10.01 on 3 and 1466 DF,  p-value: 1.57e-06
```

```r
library(stargazer)
```

```
##
## Please cite as:

##  Hlavac, Marek (2022). stargazer: Well-Formatted Regression and Summary Statistics Tables.

##  R package version 5.2.3. https://CRAN.R-project.org/package=stargazer
```

```r
stargazer(fit.reg.r, fit.reg.a, fit.reg.s,
          se= list(summary(fit.reg.r)$coef[,2],
                   summary(fit.reg.a)$coef[,2],
                   summary(fit.reg.s)$coef[,2]),
          dep.var.labels = c("Robbery", "Assault","Sexaul Assault"),
          add.lines = list(c("Standard Errors", "Standard", "Standard", "Standard")),
          dep.var.caption = "",
```

```
        no.space = T,
        header = F)
```

Table 1:

|  | Robbery | Assault | Sexaul Assault |
|---|---|---|---|
|  | (1) | (2) | (3) |
| treat | −0.010 | 0.079*** | 0.023* |
|  | (0.034) | (0.028) | (0.012) |
| post | −0.034 | 0.020 | −0.006 |
|  | (0.035) | (0.029) | (0.012) |
| treat:post | 0.045 | −0.037 | 0.039** |
|  | (0.049) | (0.041) | (0.017) |
| Constant | 0.328*** | 0.156*** | 0.011 |
|  | (0.025) | (0.020) | (0.009) |
| Standard Errors | Standard | Standard | Standard |
| Observations | 1,445 | 1,499 | 1,470 |
| $R^2$ | 0.001 | 0.007 | 0.020 |
| Adjusted $R^2$ | −0.001 | 0.005 | 0.018 |
| Residual Std. Error | 0.465 (df = 1441) | 0.396 (df = 1495) | 0.163 (df = 1466) |
| F Statistic | 0.423 (df = 3; 1441) | 3.377** (df = 3; 1495) | 10.009*** (df = 3; 1466) |

*Note:* $^{*}$p<0.1; $^{**}$p<0.05; $^{***}$p<0.01

test the parallel trends assumption with data in 2016

```r
library(ggplot2)

# robbery
ave.prepre.1.r <- with(d_matched.rb, mean(robbery[year == 2015 & treat == 1]))
ave.prepre.0.r <- with(d_matched.rb, mean(robbery[year == 2015 & treat == 0]))
se.prepre.1.r <- with(d_matched.rb, sd(robbery[year == 2015 & treat == 1])/
sqrt(sum(year == 2015 & treat == 1)))
se.prepre.0.r <- with(d_matched.rb, sd(robbery[year == 2015 & treat == 0])/
sqrt(sum(year == 2015 & treat == 0)))

# average
ave.pre.1.r <- with(r, mean(robbery[year == 2017 & treat == 1]))
ave.post.1.r <- with(r, mean(robbery[year == 2019 & treat == 1]))
ave.pre.0.r <- with(r, mean(robbery[year == 2017 & treat == 0]))
ave.post.0.r <- with(r, mean(robbery[year == 2019 & treat== 0]))

# se
se.pre.1.r <- with(r, sd(robbery[year == 2017 & treat == 1])/
sqrt(sum(year == 2017 & treat == 1)))
se.post.1.r <- with(r, sd(robbery[year == 2019 & treat == 1])/
sqrt(sum(year == 2019 & treat == 1)))
se.pre.0.r <- with(r, sd(robbery[year == 2017 & treat == 0])/
sqrt(sum(year == 2017 & treat == 0)))
se.post.0.r <- with(r, sd(robbery[year == 2019 & treat == 0])/
sqrt(sum(year == 2019 & treat == 0)))

ave.prepre.0.r
```

```
## [1] 0.4023904
```

```r
estimate.r <- c(ave.prepre.0.r,
                ave.prepre.1.r,
                ave.pre.0.r,
                ave.pre.1.r,
                ave.post.0.r,
                ave.post.1.r)
upper.r <- c(ave.prepre.0.r + 2*se.prepre.0.r,
             ave.prepre.1.r + 2*se.prepre.1.r,
             ave.pre.0.r + 2*se.pre.0.r,
             ave.pre.1.r + 2*se.pre.1.r,
             ave.post.0.r + 2*se.post.0.r,
             ave.post.1.r + 2*se.post.1.r)
lower.r <- c(ave.prepre.0.r - 2*se.prepre.0.r,
             ave.prepre.1.r - 2*se.prepre.1.r,
             ave.pre.0.r - 2*se.pre.0.r,
             ave.pre.1.r - 2*se.pre.1.r,
             ave.post.0.r - 2*se.post.0.r,
             ave.post.1.r - 2*se.post.1.r)
year = c(2015, 2015, 2017, 2017, 2019, 2019)
treat = c(0,1,0,1,0,1)
plotdata.r <- data.frame(estimate = estimate.r,
                         upper = upper.r,
                         lower = lower.r,
                         year = year,
                         treat = as.factor(treat))

v1 <- ggplot(plotdata.r,
       aes(x = year, y = estimate, color = treat)) +
  geom_point(size = 3) +
  geom_line() +
  geom_errorbar(aes(x = year, ymax = upper, ymin = lower),
                width = 0.2) +
  geom_vline(xintercept = 2018) +
  ggtitle("Robbery")
```

Oof, let's look at assault

```r
# assault
ave.prepre.1.a <- with(d_matched.as, mean(assault[year == 2015 & treat == 1]))
ave.prepre.0.a <- with(d_matched.as, mean(assault[year == 2015 & treat == 0]))
se.prepre.1.a <- with(d_matched.as, sd(assault[year == 2015 & treat == 1])/
sqrt(sum(year == 2015 & treat == 1)))
se.prepre.0.a <- with(d_matched.as, sd(assault[year == 2015 & treat == 0])/
sqrt(sum(year == 2015 & treat == 0)))

# average
ave.pre.1.a <- with(a, mean(assault[year == 2017 & treat == 1]))
ave.post.1.a <- with(a, mean(assault[year == 2019 & treat == 1]))
ave.pre.0.a <- with(a, mean(assault[year == 2017 & treat == 0]))
ave.post.0.a <- with(a, mean(assault[year == 2019 & treat== 0]))

# se
```

```r
se.pre.1.a <- with(a, sd(assault[year == 2017 & treat == 1]))/
sqrt(sum(year == 2017 & treat == 1)))
se.post.1.a <- with(a, sd(assault[year == 2019 & treat == 1]))/
sqrt(sum(year == 2019 & treat == 1)))
se.pre.0.a <- with(a, sd(assault[year == 2017 & treat == 0]))/
sqrt(sum(year == 2017 & treat == 0)))
se.post.0.a <- with(a, sd(assault[year == 2019 & treat == 0]))/
sqrt(sum(year == 2019 & treat == 0)))


estimate.a <- c(ave.prepre.0.a,
                ave.prepre.1.a,
                ave.pre.0.a,
                ave.pre.1.a,
                ave.post.0.a,
                ave.post.1.a)
upper.a <- c(ave.prepre.0.a + 2*se.prepre.0.a,
             ave.prepre.1.a + 2*se.prepre.1.a,
             ave.pre.0.a + 2*se.pre.0.a,
             ave.pre.1.a + 2*se.pre.1.a,
             ave.post.0.a + 2*se.post.0.a,
             ave.post.1.a + 2*se.post.1.a)
lower.a <- c(ave.prepre.0.a - 2*se.prepre.0.a,
             ave.prepre.1.a - 2*se.prepre.1.a,
             ave.pre.0.a - 2*se.pre.0.a,
             ave.pre.1.a - 2*se.pre.1.a,
             ave.post.0.a - 2*se.post.0.a,
             ave.post.1.a - 2*se.post.1.a)
plotdata.a <- data.frame(estimate = estimate.a,
                         upper = upper.a,
                         lower = lower.a,
                         year = year,
                         treat = as.factor(treat))

v2 <- ggplot(plotdata.a,
      aes(x = year, y = estimate, color = treat)) +
  geom_point(size = 3) +
  geom_line() +
  geom_errorbar(aes(x = year, ymax = upper, ymin = lower),
                width = 0.2) +
  geom_vline(xintercept = 2018) +
  ggtitle("Assault")
```

poo, let's look at sexual assault now

```r
# sexual assault
ave.prepre.1.s <- with(d_matched.sa, mean(sexual.assault[year == 2015 & treat == 1]))
ave.prepre.0.s <- with(d_matched.sa, mean(sexual.assault[year == 2015 & treat == 0]))
se.prepre.1.s <- with(d_matched.sa, sd(sexual.assault[year == 2015 & treat == 1]))/
sqrt(sum(year == 2015 & treat == 1)))
se.prepre.0.s <- with(d_matched.sa, sd(sexual.assault[year == 2015 & treat == 0]))/
sqrt(sum(year == 2015 & treat == 0)))

# average
```

```r
ave.pre.1.s <- with(s, mean(sexual.assault[year == 2017 & treat == 1]))
ave.post.1.s <- with(s, mean(sexual.assault[year == 2019 & treat == 1]))
ave.pre.0.s <- with(s, mean(sexual.assault[year == 2017 & treat == 0]))
ave.post.0.s <- with(s, mean(sexual.assault[year == 2019 & treat== 0]))

# se
se.pre.1.s <- with(s, sd(sexual.assault[year == 2017 & treat == 1])/
sqrt(sum(year == 2017 & treat == 1)))
se.post.1.s <- with(s, sd(sexual.assault[year == 2019 & treat == 1])/
sqrt(sum(year == 2019 & treat == 1)))
se.pre.0.s <- with(s, sd(sexual.assault[year == 2017 & treat == 0])/
sqrt(sum(year == 2017 & treat == 0)))
se.post.0.s <- with(s, sd(sexual.assault[year == 2019 & treat == 0])/
sqrt(sum(year == 2019 & treat == 0)))


estimate.s <- c(ave.prepre.0.s,
                ave.prepre.1.s,
                ave.pre.0.s,
                ave.pre.1.s,
                ave.post.0.s,
                ave.post.1.s)
upper.s <- c(ave.prepre.0.s + 2*se.prepre.0.s,
             ave.prepre.1.s + 2*se.prepre.1.s,
             ave.pre.0.s + 2*se.pre.0.s,
             ave.pre.1.s + 2*se.pre.1.s,
             ave.post.0.s + 2*se.post.0.s,
             ave.post.1.s + 2*se.post.1.s)
lower.s <- c(ave.prepre.0.s - 2*se.prepre.0.s,
             ave.prepre.1.s - 2*se.prepre.1.s,
             ave.pre.0.s - 2*se.pre.0.s,
             ave.pre.1.s - 2*se.pre.1.s,
             ave.post.0.s - 2*se.post.0.s,
             ave.post.1.s - 2*se.post.1.s)
plotdata.s <- data.frame(estimate = estimate.s,
                         upper = upper.s,
                         lower = lower.s,
                         year = year,
                         treat = as.factor(treat))

v3 <- ggplot(plotdata.s,
       aes(x = year, y = estimate, color = treat)) +
  geom_point(size = 3) +
  geom_line() +
  geom_errorbar(aes(x = year, ymax = upper, ymin = lower),
                width = 0.2) +
  geom_vline(xintercept = 2018) +
  ggtitle("Sexual Assault")

library(gridExtra)
grid.arrange(v1, v2, v3, ncol=2)
```
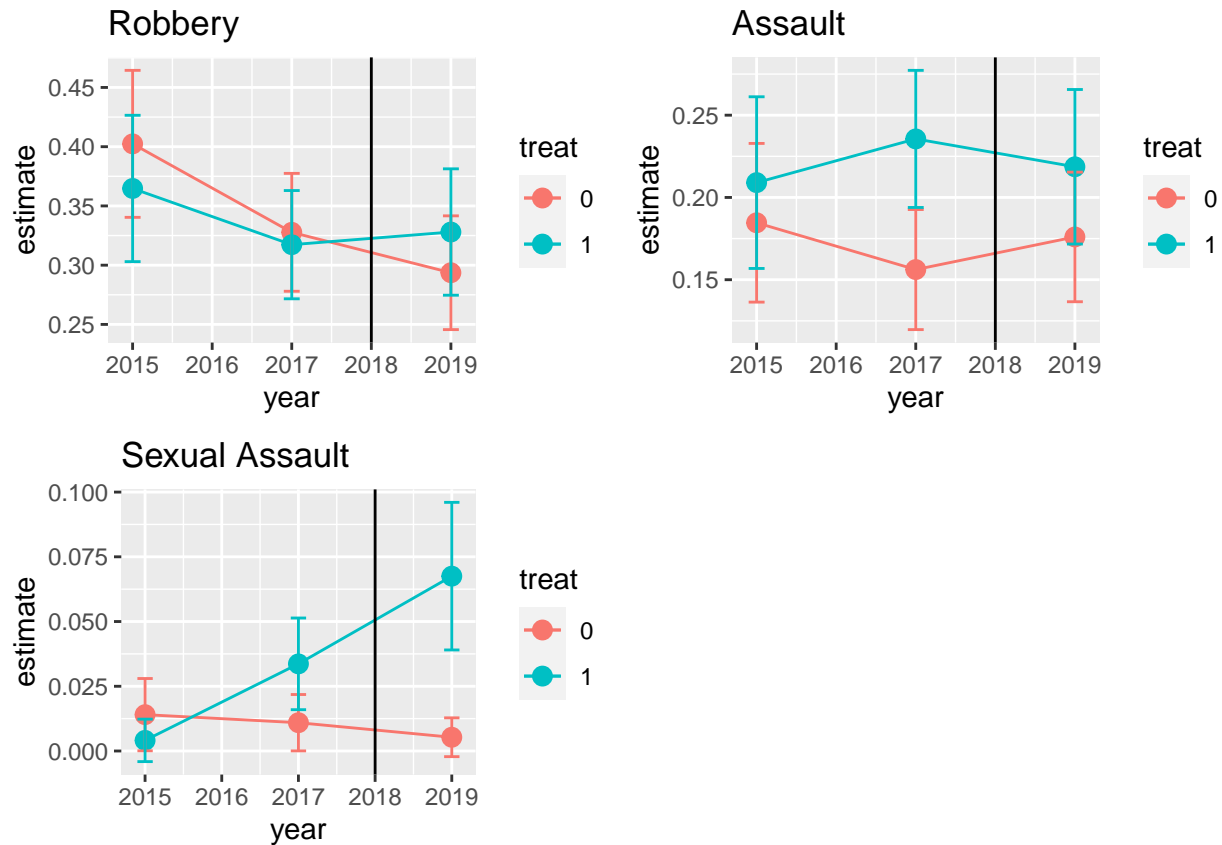
## Robbery

## Assault

## Sexual Assault

now run placebo test using 2016 as the placebo year

```r
# robbery
placebo.r <- subset(d_matched.rb, (year > 2014) & (year < 2018))
placebo.r$post <- ifelse(placebo.r$year > 2016, 1, 0)

dum1 <- lm(robbery ~ treat * post, data = placebo.r)

# assault
placebo.a <- subset(d_matched.as, (year > 2014) & (year < 2018))
placebo.a$post <- ifelse(placebo.a$year > 2016, 1, 0)

dum2 <- lm(assault ~ treat * post, data = placebo.a)

# sexual assault
placebo.s <- subset(d_matched.sa, (year > 2014) & (year < 2018))
placebo.s$post <- ifelse(placebo.s$year > 2016, 1, 0)

dum3 <- lm(sexual.assault ~ treat * post, data = placebo.s)


stargazer(dum1, dum2, dum3,
        se= list(summary(dum1)$coef[,2],
                summary(dum2)$coef[,2],
                summary(dum3)$coef[,2]),
        dep.var.labels = c("Robbery Dummy", "Assault Dummy","Sexaul Assault Dummy"),
        add.lines = list(c("Standard Errors", "Standard", "Standard", "Standard")),
```

```
        dep.var.caption = "",
        no.space = T,
        header = F)
```

Table 2:

|  | Robbery Dummy | Assault Dummy | Sexaul Assault Dummy |
|---|---|---|---|
|  | (1) | (2) | (3) |
| treat | −0.042 | 0.002 | 0.006 |
|  | (0.028) | (0.023) | (0.008) |
| post | −0.072** | −0.023 | −0.001 |
|  | (0.032) | (0.025) | (0.009) |
| treat:post | 0.031 | 0.077** | 0.016 |
|  | (0.044) | (0.036) | (0.012) |
| Constant | 0.399*** | 0.179*** | 0.012** |
|  | (0.020) | (0.016) | (0.005) |
| Standard Errors | Standard | Standard | Standard |
| Observations | 1,939 | 1,990 | 2,018 |
| $R^2$ | 0.004 | 0.005 | 0.004 |
| Adjusted $R^2$ | 0.003 | 0.003 | 0.002 |
| Residual Std. Error | 0.478 (df = 1935) | 0.389 (df = 1986) | 0.134 (df = 2014) |
| F Statistic | 2.901** (df = 3; 1935) | 3.115** (df = 3; 1986) | 2.609** (df = 3; 2014) |

*Note:* $^{*}$p<0.1; $^{**}$p<0.05; $^{***}$p<0.01

regression by race

```
library(dplyr)
```

```
##
## Attaching package: 'dplyr'

## The following object is masked from 'package:gridExtra':
##
##     combine

## The following objects are masked from 'package:stats':
##
##     filter, lag

## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
```

```
load("sa_no_race.RData")

# subset data
nr <- subset(d_matched.sa.nr, (year > 2016) & (year < 2020))
# remove treatment year
nr <- subset(nr, year != 2018)

# use regression to estimate difference in difference
nr$post <- ifelse(nr$year > 2018, 1, 0)

# only get the first 5 race categories
```

```r
nr <- subset(nr, race<6)

nr.w <- subset(nr, race==1)
nr.b <- subset(nr, race==2)
nr.na <- subset(nr, race==3)
nr.as <- subset(nr, race==4)
nr.pc <- subset(nr, race==5)

nr.w.r <- lm(sexual.assault ~ treat * post, data = nr.w)
nr.b.r <- lm(sexual.assault ~ treat * post, data = nr.b)
nr.na.r <- lm(sexual.assault ~ treat * post, data = nr.na)
nr.as.r <- lm(sexual.assault ~ treat * post, data = nr.as)
nr.pc.r <- lm(sexual.assault ~ treat * post, data = nr.pc)

stargazer(nr.w.r, nr.b.r, nr.na.r,
          se= list(summary(nr.w.r)$coef[,2],
                   summary(nr.b.r)$coef[,2],
                   summary(nr.na.r)$coef[,2]),
          dep.var.labels = c("White", "Black","Native american alaskan"),
          add.lines = list(c("Standard Errors", "Standard", "Standard", "Standard")),
          dep.var.caption = "",
          no.space = T,
          header = F)
```

Table 3:

|  | White | | |
|---|---|---|---|
|  | (1) | (2) | (3) |
| treat | 0.008 | $-0.037$ | $-0.000$ |
|  | (0.015) | (0.029) | (0.115) |
| post | $-0.021$ | $-0.037$ | 0.083 |
|  | (0.016) | (0.029) | (0.116) |
| treat:post | 0.057** | 0.054 | 0.417** |
|  | (0.022) | (0.041) | (0.182) |
| Constant | 0.029*** | 0.052** | 0.000 |
|  | (0.011) | (0.022) | (0.093) |
| Standard Errors | Standard | Standard | Standard |
| Observations | 1,070 | 254 | 36 |
| $R^2$ | 0.015 | 0.008 | 0.303 |
| Adjusted $R^2$ | 0.012 | $-0.004$ | 0.238 |
| Residual Std. Error | 0.182 (df = 1066) | 0.164 (df = 250) | 0.245 (df = 32) |
| F Statistic | 5.284*** (df = 3; 1066) | 0.699 (df = 3; 250) | 4.638*** (df = 3; 32) |

*Note:* $^*$p<0.1; $^{**}$p<0.05; $^{***}$p<0.01

```r
# clearly, American Indian, Alaska native only is the only statistically significant field. It is also

# let's run the dummy variable

native_res <- subset(nr, race==3) # 39 data points

native <- subset(d_matched.sa.nr, race==3) # 286 data points
```

```
placebo.s.native <- subset(native,
                           (year > 2014) & (year < 2018))

placebo.s.native$post <- ifelse(placebo.s.native$year > 2016, 1, 0)

summary(lm(sexual.assault ~ treat * post, data = placebo.s.native))
```

Call: lm(formula = sexual.assault ~ treat * post, data = placebo.s.native)

Residuals: Min 1Q Median 3Q Max 0 0 0 0 0

Coefficients: Estimate Std. Error t value Pr(>|t|) (Intercept) 0 0 NaN NaN treat 0 0 NaN NaN post 0 0 NaN NaN treat:post 0 0 NaN NaN

Residual standard error: 0 on 50 degrees of freedom Multiple R-squared: NaN, Adjusted R-squared: NaN F-statistic: NaN on 3 and 50 DF, p-value: NA

```
# this is because no one reported sexaul assualt
```

plot before and after for sexual assault

```
library(plyr)
```

```
## --------------------------------------------------------------------------------
## You have loaded plyr after dplyr - this is likely to cause problems.
## If you need functions from both plyr and dplyr, please load plyr first, then dplyr:
## library(plyr); library(dplyr)
## --------------------------------------------------------------------------------
##
## Attaching package: 'plyr'

## The following objects are masked from 'package:dplyr':
##
##     arrange, count, desc, failwith, id, mutate, rename, summarise,
##     summarize
```

```
s <- subset(s, race < 6)


s <- aggregate(s$sexual.assault,
               by = list(YearQ = s$year.q,
                         race = s$race,
                         treat = s$treat,
                         post = s$post),
               FUN = mean)

s$post <- as.character(s$post)
# map from integer to string
raceNames <- list(
  "1"="White Only",
  "2"="Black Only",
  "3"="Native American Alaskan Native",
  "4"="Asian",
  "5"="Pacific Islander"
)
```
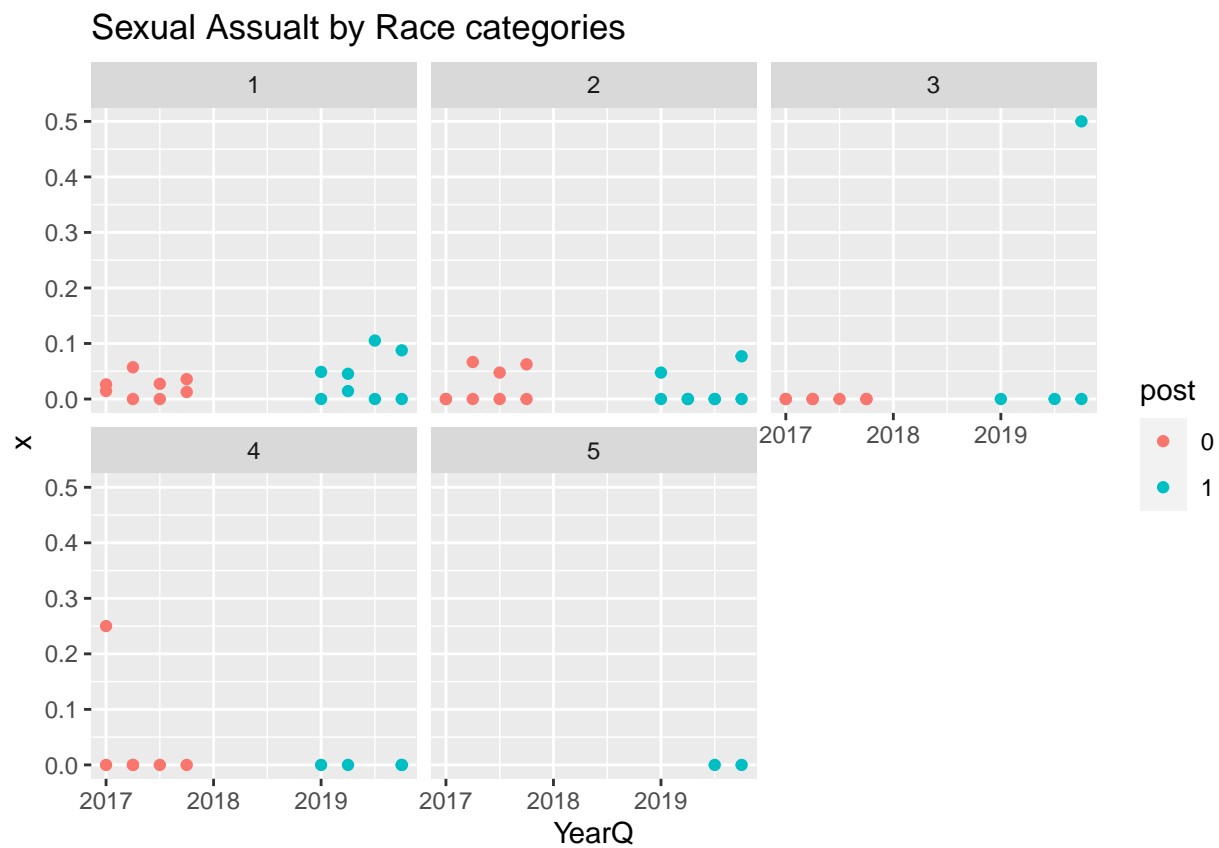
10

```
s$race <- as.integer(s$race)

labeler <- function (variable, value){
  return (raceNames[value])
}

plot1 <- ggplot() +
  geom_point(data = s, aes(x = YearQ, y = x, color = post)) +
  facet_wrap(s$race) +
  ggtitle("Sexual Assualt by Race categories")

plot1
```



Sexual Assualt by Race categories

## Robustness testing

This code replicates cleaning, matching, and diff-in-diff with revised categories for identifying treatment, for the purpose of robustness matching. Outcome variable is sexual assault

# Data cleaning

First, clean data and make treatment, outcome variables

```r
# first, load the data
library(stringi)

load("~/Documents/2021-2022/Spring2022/17.803/project/test3.rda")

# columns of interest:
c.interest <- c(
  "YEARQ",
  "YEAR", # year and quarter of interview
  "V2014", # tenure, owned house/cash rent/no cash rent
  "V2016", # land use, urban/rural/residue
  "V2020", # type of living quarters
  "V2024", # number of housing units in structure
  "V2025", # direct outside access
  "V2025A", # gated or walled community
  "V2025B", # building with restricted access
  "V2026", # household income
  "V2033", # principle person age
  "V2036", # principle person sex
  "V2038", # principle person educational attainment
  "V2040A", # principle person race recode
  "V2074", # operate business from address
  "V2075", # sign on premises for business
  "V2122", # family structure code
  "V3033", # times moved in the last 5 years
  # started since 2017: -------------
  "V3085", # gender identity at birth
  "V3086", # current gender identity
  # end of 2017 only block --------------
  "V4529" # Type of crime code new
)

intersect <- intersect(names(da38136.0003), c.interest)
s <- da38136.0003[,intersect]
s <- s[!is.na(s$YEARQ),]

s$YEARQ <- as.character(s$YEARQ)

# first quarter
s$YEARQ <- gsub("\\b.1\\b", "", s$YEARQ)

# second quarter
s$YEARQ <- gsub("\\b.2\\b", ".25", s$YEARQ)

# third quarter
```

```r
s$YEARQ <- gsub("\\b.3\\b", ".50", s$YEARQ)

# fourth quarter
s$YEARQ <- gsub("\\b.4\\b", ".75", s$YEARQ)

# back to numeric:
s$YEARQ <- as.numeric(s$YEARQ)

# make outcome variable: sexual assault
sexual.assault <- c(
  "\\<(01)\\>",
  "\\<(02)\\>",
  "\\<(03)\\>",
  "\\<(04)\\>",
  "\\<(15)\\>",
  "\\<(16)\\>",
  "\\<(18)\\>",
  "\\<(19)\\>"
)

s$sexual.assault <- grepl(paste(sexual.assault, collapse = "|"), s$V4529)
# convert from boolean to integer
s$sexual.assault <- as.integer(s$sexual.assault)

# start mapping indicator variables

# living quarters
living.quarters <- c(
  "\\<(02)\\>",
  "\\<(03)\\>",
  "\\<(04)\\>",
  "\\<(05)\\>",
  "\\<(06)\\>",
  "\\<(07)\\>",
  "\\<(08)\\>",
  "\\<(09)\\>",
  "\\<(10)\\>"
)

number.units <- c(
  "\\<(01)\\>",
  # "\\<(02)\\>",
  "\\<(98)\\>",
  "\\<(99)\\>"
)

income <- c(
  "\\<(13)\\>",
  "\\<(14)\\>",
  "\\<(15)\\>",
  "\\<(16)\\>",
  "\\<(17)\\>",
  "\\<(18)\\>",
```

```r
  "\\<(99)\\>"
)

# trans indicator
s$V3085 <- as.character(s$V3085)
s$V3086 <- as.character(s$V3086)

trans.discount <- c(
  "\\b(3)\\b",
  "\\b(8)\\b",
  "\\b(-1)\\b",
  "\\b(9)\\b"
)

s$trans <-
  (!is.na(s$V3085) & !is.na(s$V3086)) &
  (s$V3085 != s$V3086) &
  !grepl(paste(trans.discount,
               collapse = "|"), s$V3085) &
  !grepl(paste(trans.discount,
               collapse = "|"), s$V3086)

s$treat <-
  ((s$V2014=="(2) Rented for cash" | s$V2014=="(3) No cash rent") &
  grepl(paste(living.quarters, collapse = "|"), s$V2020) & # 49771
  !grepl(paste(number.units, collapse = "|"), s$V2024) & # 49771
  s$V2025A == "(2) No" & # not gated/walled community 42813
  s$V2025B == "(2) No" & # no restricted access 39760
  (s$V2033 > 15 & s$V2033 < 51) & # age range 36436
  (as.character(s$V2036)=="(2) Female"| s$trans == TRUE) & # either female or transgender 25923
  grepl("Lone", s$V2122) &  # lone male/female
  s$V3033>2 & # moved more than twice in past 5 years 13613
  !grepl(paste(income, collapse = "|"), s$V2026)) |  # income less than 50000
  (s$V2074 == "(01) Yes" & s$V2075 == "(02) No" ) # or operate business but no sign)

s <- s[, c(
  "YEARQ",
  "YEAR",
  "treat",
  "sexual.assault",
  "V2038",
  "V2016",
  "V2040A"
)]
save(s, file="Robust.RData")
```

## Matching

Now, load s back and match

```r
load("Robust.RData")

library(Matching)
```

```
## Loading required package: MASS

## ##
## ##  Matching (Version 4.9-11, Build Date: 2021-10-18)
## ##  See http://sekhon.berkeley.edu/matching for additional documentation.
## ##  Please cite software as:
## ##   Jasjeet S. Sekhon. 2011. ``Multivariate and Propensity Score Matching
## ##   Software with Automated Balance Optimization: The Matching package for R.''
## ##   Journal of Statistical Software, 42(7): 1-52.
## ##
```

```r
library(ebal)
```

```
## ##
## ## ebal Package: Implements Entropy Balancing.

## ## See http://www.stanford.edu/~jhain/ for additional information.
```

```r
library(ggplot2)
library(dplyr)
```

```
##
## Attaching package: 'dplyr'

## The following object is masked from 'package:MASS':
##
##     select

## The following objects are masked from 'package:stats':
##
##     filter, lag

## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
```

```r
library(xtable)
library(tidyr)
library(foreign)

s.no.na <- s[!is.na(s$V2040A),]
s.no.na <- s.no.na[!is.na(s.no.na$V2016),]
s.no.na <- s.no.na[!is.na(s.no.na$V2038),]
s.no.na <- s.no.na[!is.na(s.no.na$treat),]

# convert categorical to numerical
s.no.na$V2038 <- unclass(s.no.na$V2038)
s.no.na$V2016 <- unclass(s.no.na$V2016)
s.no.na$V2040A <- unclass(s.no.na$V2040A)
s.no.na$treat <- as.integer(s.no.na$treat )

colnames(s.no.na) <- c(
  "year.q",
  "year",
  "treat",
  "sexual.assault",
  "education",
  "urban.rural",
```

```
    "race"
)

# match
SW.formula.sa <- formula(treat ~ education + urban.rural + race)

varnames.sa <- c("education", "urban.rural", "race")

match.sa <- Match(Y = s.no.na$sexual.assault,
                  Tr = s.no.na$treat,
                  X = s.no.na[,varnames.sa],
                  Weight = 2,
                  ties=FALSE)

matched.sa <- MatchBalance(SW.formula.sa,
                           match.out = match.sa,
                           data = s.no.na,
                           print.level = 0)

d_matched.sa <- s.no.na[c(match.sa$index.treated, match.sa$index.control),]
try <- subset(d_matched.sa, sexual.assault==1)
save(d_matched.sa, file="RobustAfterMatching.RData")
```

## diff-in-diff

```
load("RobustAfterMatching.RData")

# at this point, we realize that we are in troulbe
# because now we have no outcome units that are 1

# backtrack strategy: go back to pre-matching, and instead do regression
# with covariates

load("Robust.RData")

s.no.na <- s[!is.na(s$V2040A),]
s.no.na <- s.no.na[!is.na(s.no.na$V2016),]
s.no.na <- s.no.na[!is.na(s.no.na$V2038),]

# convert categorical to numerical
s.no.na$V2038 <- unclass(s.no.na$V2038)
s.no.na$V2016 <- unclass(s.no.na$V2016)
s.no.na$V2040A <- unclass(s.no.na$V2040A)

colnames(s.no.na) <- c(
  "year.q",
  "year",
  "treat",
  "sexual.assault",
  "education",
  "urban.rural",
  "race"
```

```
)

s <- subset(s.no.na, (year > 2016) & (year < 2020))
s <- subset(s, year != 2018)
s$post <- ifelse(s$year > 2018, 1, 0)

fit.reg.s <- lm(sexual.assault ~ treat * post +
                    education +
                    urban.rural +
                    race, data = s)
summary(fit.reg.s)
```

Call: lm(formula = sexual.assault ~ treat * post + education + urban.rural + race, data = s)

Residuals: Min 1Q Median 3Q Max -0.02489 -0.01803 -0.01709 -0.01616 0.99195

Coefficients: Estimate Std. Error t value Pr(>|t|)
(Intercept) 3.115e-02 5.886e-03 5.292 1.22e-07 * **treatTRUE -1.141e-02 5.718e-02 -0.200 0.8418
post 1.250e-03 1.763e-03 0.709 0.4781
education -3.116e-04 2.196e-04 -1.419 0.1560
urban.rural -6.829e-03 2.475e-03 -2.759 0.0058** race -3.696e-04 5.749e-04 -0.643 0.5203
treatTRUE:post -2.222e-06 9.333e-02 0.000 1.0000
— Signif. codes: 0 '*** *0.001 ** *0.01* *'* 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1278 on 21073 degrees of freedom Multiple R-squared: 0.0004727, Adjusted R-squared: 0.0001881 F-statistic: 1.661 on 6 and 21073 DF, p-value: 0.1262

```
library(stargazer)
```

```
##
## Please cite as:
```

```
##  Hlavac, Marek (2022). stargazer: Well-Formatted Regression and Summary Statistics Tables.
```

```
##  R package version 5.2.3. https://CRAN.R-project.org/package=stargazer
```

```
stargazer(fit.reg.s,
          se= list(summary(fit.reg.s)$coef[,2]),
          dep.var.labels = c("Sexaul Assault"),
          add.lines = list(c("Standard Errors", "Standard")),
          dep.var.caption = "",
          no.space = T,
          header = F)
```

Table 1:

|  | Sexaul Assault |
|---|---|
| treat | −0.011 |
|  | (0.057) |
| post | 0.001 |
|  | (0.002) |
| education | −0.0003 |
|  | (0.0002) |
| urban.rural | −0.007*** |
|  | (0.002) |
| race | −0.0004 |
|  | (0.001) |
| treatTRUE:post | −0.00000 |
|  | (0.093) |
| Constant | 0.031*** |
|  | (0.006) |
| Standard Errors | Standard |
| Observations | 21,080 |
| R$^2$ | 0.0005 |
| Adjusted R$^2$ | 0.0002 |
| Residual Std. Error | 0.128 (df = 21073) |
| F Statistic | 1.661 (df = 6; 21073) |

*Note:* $^*$p<0.1; $^{**}$p<0.05; $^{***}$p<0.01