

Final Report: U.S. Gun Violence Trends

By Jonathan Shuai

[Problem Statement](#)

[Data Wrangling](#)

[I. Summary of the Data Used](#)

[II. Cleaning and Consolidating the Data](#)

[III. Null Values, Outliers, and Provisions Data](#)

[Null Values](#)

[Outliers](#)

[Provisions](#)

[Exploratory Data Analysis and Hypothesis Testing](#)

[I. Provision Efficiency](#)

[Single Provisions](#)

[Pairs of Provisions](#)

[II. Testing for Structural Breaks](#)

[Predictive Modeling](#)

[I. Problem Structure and Label](#)

[II. Monthly Modeling](#)

[Pipeline and Modeling](#)

[Results and Challenges](#)

[III. Annual Modeling](#)

[Pipeline and Modeling](#)

[Results and Challenges](#)

[Conclusion](#)

Problem Statement

The United States has notoriously high gun violence and mass shootings incidents, and the problem seems to be getting worse over time. Often people say that gun control laws wouldn't affect gun violence rates because gun homicides are committed by criminals who don't follow laws. However, I wanted to see if analyzing the data would give help to confirm or reject this claim.

The primary goal of this project is to examine gun violence trends in the United States and find correlations with gun control laws. Here are some of the questions we will be answering:

- How have gun violence rates changed over the years?
Which states have seen the biggest increases in gun violence? the lowest?
- How does gun control in states with low gun violence rates compare with gun control in states with high gun violence rates?
- Which gun control laws have the most correlation with reduced gun violence?
- Are there categories of gun control laws that perform better than others? (e.g. background checks vs. banning assault weapon?)

We are also interested to see if we can find any relationship between gun violence and features such as income, substance abuse, and other crime. Finally, we will be creating a model to help predict whether gun violence will increase in the next year for each state.

Through this project, I hope to find useful insights on gun control provisions that may help policymakers make better decisions in interest of reducing gun violence. Using data to analyze laws allows us to justify policies with evidence and understand which laws work and which ones don't. Furthermore, by predicting gun violence increases in each state, we can take countermeasures to prevent them such as budgeting for law enforcement or implementing new gun control policies.

Data Wrangling

I. Summary of the Data Used

The data used in this project comes from many different sources. The primary goal of the data wrangling was to create CSV files that consolidate all of the useful features to be used during the visualization and modeling stages.

Here are the following datasets used listed alphabetically. Each entry includes the name of the file in the './data/raw' directory, the source, and descriptions of each dataset:

[National Institute on Alcoholism and Alcohol Abuse](#)

Alcohol consumption per capita for each state from 1977-2016.

./data/raw/alcohol.csv

[GunPolicy.org](#)

Annual gun homicides data per state from 2000 to 2013.

./data/raw/annual_gun_deaths.csv

[Gun Violence Archive](#)

Gun violence incidents from January 2014 to March 2018. This dataset was downloaded from Kaggle, although it was originally scraped from the Gun Violence Archive.

**Note: I added one entry for the Las Vegas shooting in October 2017, as it appeared to be missing.*

./data/raw/incidents.csv.gz

[Disaster Center](#)

Annual crime factors from 2010 to 2016. The Disaster Center collected the data from the FBI UCS Annual Crime Reports.

./data/raw/crime.csv

[US Census Bureau](#)

Two different population CSVs were downloaded from the US Census Bureau, and merged together in a simpler file ./data/raw/population.csv via a Python script merge-populations.py. This merged CSV contains annual populations for each state from 2000 to 2017.

./data/raw/population_2000_2010.csv

./data/raw/population_2010_2017.csv

./data/raw/population.csv

[Kaggle](#)

Gun control provisions as annual entries for each state from 1991 to 2017. This dataset was generated from several sources, including Thomson Reuters Westlaw legislative database and data from Everytown for Gun Safety and Legal Science, LLC. Each gun provision is encoded as a shortened codename. The details about each provision and its shortened name can be found in ./data/raw/codebook.xlsx

./data/raw/provisions.csv

[270 to Win](#)

Election results from 2000 to 2016 for each state. This dataset simply lists the percentage of votes each part received in each state.

./data/raw/election_results.csv

[Bureau of Economic Analysis](#)

Personal income data from 2009-2017, listed annually and for each state.

`./data/raw/income.csv`

[Bureau of Alcohol, Tobacco, and Firearms](#)

Annual gun registration data for each state from 2011-2017.

`./data/raw/registrations.csv`

[SAMHSA](#)

Substance use for each state from survey results from 2012-2016. 6 features were selected from the survey results. They are:

marijuana - Use of marijuana in the last year

cocaine - Use of cocaine in the last year

tobacco - Use of tobacco within the last month

alcohol - Alcohol abuse or dependency

mental - Any mental illness

depression - Serious depression episode within the last year

More details of the criteria of these features can be found from the SAMHSA website.

`./data/raw/substances.csv`

II. Cleaning and Consolidating the Data

Each dataset has a different range of years for which it is available. Here are the important things to note about the availability of the data:

- **Daily gun homicide data is available for 2014-2017**, whereas only **annual gun homicide data is available for 2000-2013**.
- **All data is available for at least 2013-2016**, allowing us to use every feature in a model predicting gun violence for 2014-2017 (we make predictions for at least one year into the future).

To make the data easier to work with, I decided to combine the data into several CSV files, each structured to make specific tasks easier later on. Here is a summary of these CSV files:

`./data/cleaned/annual.csv`

This CSV contains all of the features* for each state from 2000 - 2017, with entries where feature was not available as the only null values. This will be the primary CSV used for visualizations, as it is easy to interpret and manipulate.

`./data/cleaned/feature.csv`

This CSV contains monthly homicide rates for each state from 2014-2017. It is organized as gun homicide observations for each month and state, paired with all of the annual features* from the previous year. This will be the CSV we use for modeling,

as it has no null values. We will be using observations in 2014-2016 for training to predict monthly gun homicide rate changes for each state in 2017.

./data/cleaned/by_date_total.csv

A CSV that contains the number of daily deaths for each state (as columns) from January 2014 to March 2018. I decided to make this CSV because it is the only CSV that will contain daily homicide rates. These daily homicide rates will be resampled and used for visualizations as well as time series analysis.

./data/cleaned/location.csv

A CSV that contains latitude and longitude coordinates for each incident. It also includes location information for each incident. Since we have a lot of missing values for these features and we don't care about individual incidents for our modeling, I decided to put them in a separate CSV file for a few quick visualizations later.

* For the annual and feature files, provisions data was excluded. The provisions data will be added accordingly during visualization and modeling. The reasoning is explained below.

III. Null Values, Outliers, and Provisions Data

Null Values

There are several null values to be noted among the datasets:

- **'Suppressed' values for annual gun homicide totals for states with very low annual gun homicides.** There are 6 states with these suppressed values (Hawaii, New Hampshire, North Dakota, South Dakota, Vermont, and Wyoming). According to the source GunPolicy.org, these values are cases with fewer than 10 homicides. Since these states have very low gun violence incidents, I imputed them with the mean during visualization. During modeling, states with very low average gun homicide rates are excluded, including these 6 states.
- **Missing values for District of Columbia for provisions data.** There are no entries for the District of Columbia for gun control provisions. This is unfortunate as the District of Columbia is an outlier for high gun violence rates. Ultimately, the District of Columbia will have to be dropped from the analysis, as provisions data is extremely important to this project.
- **Missing latitude and longitude values for around 40% of incidents.** There seems to be more missing values the more recent the incident is. However, this is not a big problem, as these values will only be used for an fun visualization.

Outliers

In our data, the notable outliers are states like Hawaii, which has an extremely low gun violence rate, and the District of Columbia, which has an extremely high gun violence rate. These

outliers are useful in analysis and is part of our exploration on the factors that affect gun violence rates. Deciding whether to remove these outliers is an important for the modeling phase, and will be decided when tweaking the features before training.

Provisions

Furthermore, there are so many provisions (133 in total), that I decided to add them later as needed for visualization and modeling. The provisions data is not limited by availability; it contains years from 1991 to 2017, so it will be easy to selectively choose the provisions to use later. This will give us more flexibility in visualizations and modeling, and make the data cleaner overall.

Exploratory Data Analysis and Hypothesis Testing

I. Provision Efficiency

Single Provisions

Under the “Analyzing the Effects of Provisions” portion of my project, I attempt to find gun control provisions that correlate most with reduced gun violence. I examine the differences in gun violence for states that have a certain gun control provision in effect compared with states that don’t have that provision for a given year.

In order to control for sample size and outliers, I used a minimum sample size threshold and discarded provisions where the number of samples in either group was below the threshold. I created bar chart visualizations, but I wanted to conduct a t-test to confirm the differences seen in the bar chart.

The chart below shows the p-values:

Provision	Gun Violence Rates in States Without Provision	Gun Violence Rates in States With Provision	Difference	P-value
nosyg	6.024715	3.466635	2.558080	0.000893
immunity	5.629389	3.030228	2.599161	0.001455
violentpartial	5.359422	3.441462	1.917959	0.024187
statechecksh	5.336470	3.490234	1.846237	0.030354
mcdvdating	5.349055	3.574407	1.774649	0.034754
cap14	5.286320	3.596802	1.689518	0.048442

Here's a brief description of each provision codename:

nosyg - Use of deadly force is not allowed to be a first resort in public. This is sometimes referred to as a "stand your ground" law.

immunity - No law provides blanket immunity to gun manufacturers or prohibits state or local lawsuits against gun manufacturers.

violentpartial - Firearm possession is prohibited for people who have committed a violent misdemeanor punishable by more than one year of imprisonment.

mcdvdating - All people convicted of a misdemeanor crime of domestic violence are prohibited from possessing firearms.

statechecksh - State conducts separate background checks, beyond the federal NICS check, for handguns.

cap14 - Criminal liability for negligent storage applies to access by children less than 14 years old.

For a better explanation of what each provision means, the codenames can be looked up in the [./data/raw/codebook.xlsx](#)

Pairs of Provisions

I also wanted to see whether a pair of provisions would have a significant correlation with reduced gun violence.

Once again, a t-test is used to get a p-value for each of the differences:

Provision	Gun Violence Rates in States Without Provision	Gun Violence Rates in States With Provision	Difference	P-value
age21handgunsale permith	5.235497	2.698168	2.537328	0.015916
statechecks universalpermith	5.330041	2.913307	2.416734	0.021859
statechecks universalpermit	5.187620	2.913307	2.274313	0.022568
age21handgunsale universalpermith	5.336773	2.982114	2.354659	0.026356

Here are the descriptions of the provisions above:

statechecks - State conducts separate background checks, beyond the federal NICS check, for all firearm.

universalpermith - Background checks conducted through permit requirement for all handgun sales (or universal background checks).

universalpermit - Background checks conducted through permit requirement for all firearm sales (or universal background checks).

age21handgunsale - Purchase of handguns from licensed dealers and private sellers restricted to age 21 and older.

Note: The large number of provisions tested can lead to a high false positive rate. This is to be considered when looking at this analysis, and I'm presenting the findings here to show the results of looking at these features.

II. Testing for Structural Breaks

We looked at when provisions were either added or removed by a state to see if we could see any differences in gun violence after those changes were made. In order to quantify this change, we did a Chow test for structural break, an f-statistic defined as:

$$\frac{((RSS_{pooled} - (RSS_{before} + RSS_{after})) / k)}{(RSS_{before} + RSS_{after}) / (N_{before} + N_{after} - 2k)}$$

where:

RSS_{pooled} - The residual sum of squares from regression on the entire time series

RSS_{before} - The residual sum of squares from regression on the time series before the structural break

RSS_{after} - The residual sum of squares from regression on the time series after the structural break

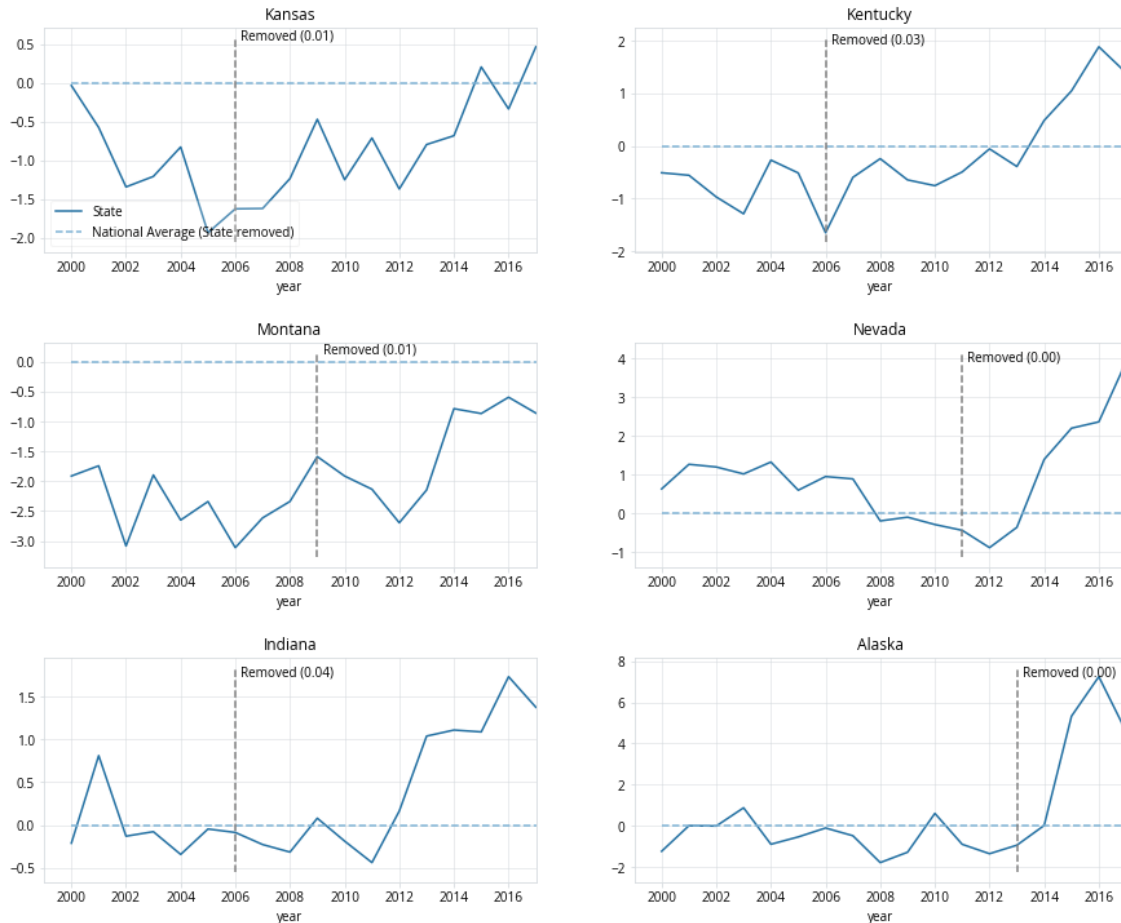
N_{before} - Number of samples in the time series before the structural break

N_{after} - Number of samples in the series after the structural break

k - Number of parameters estimated during the regression

To visualize this test, each time series is shown with a grey dotted line which indicates a change in provision. Next to it in parentheses is the p-value of the Chow test at that point. For conciseness, only one examples of these tests are shown (more can be found in [visualization.ipynb](#)):

Effects of 'nosyg' on Gun Violence



This test gives us some confidence in making judgements about observed changes in gun violence rates after a provision is changed. However, there is a limited sample size for the Chow test, and a large number of factors that could affect gun violence. The structural break test results are to be interpreted with this in mind.

Predictive Modeling

I. Problem Structure and Label

In order to use the data in a predictive model to give insight into gun violence trends, we need to define the label that we want to predict. I decided that we should set the problem up as

a classification problem rather than a regression problem because predicting the exact number of deaths would be infeasible.

I set the problem up as a binary classification with the following structure:

- Each row in our data represents an observation: features for a state in a given time period
- For each row, a 0/1 label is defined as whether **the number of gun homicide will increase from this period by more than 20% in the next time period.**
- The rows are sorted by date, and predictions will be made for each time period using **all data available** prior to that time period.

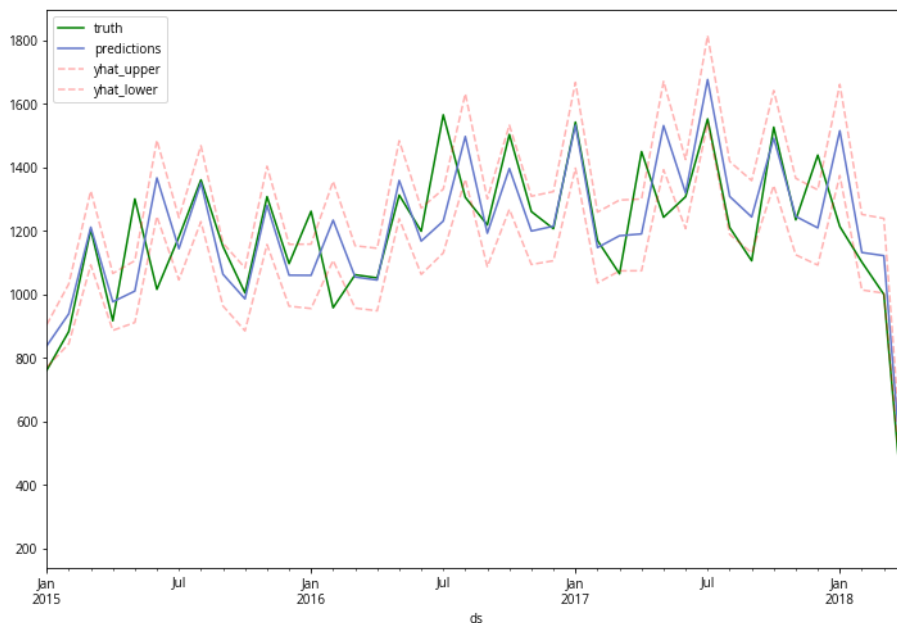
Example structure for monthly time period:

state	time_period	features...	label
Alabama	2014-02	...	0
Alaska	2014-02	...	0
...
Alabama	2016-12	...	1

II. Monthly Modeling

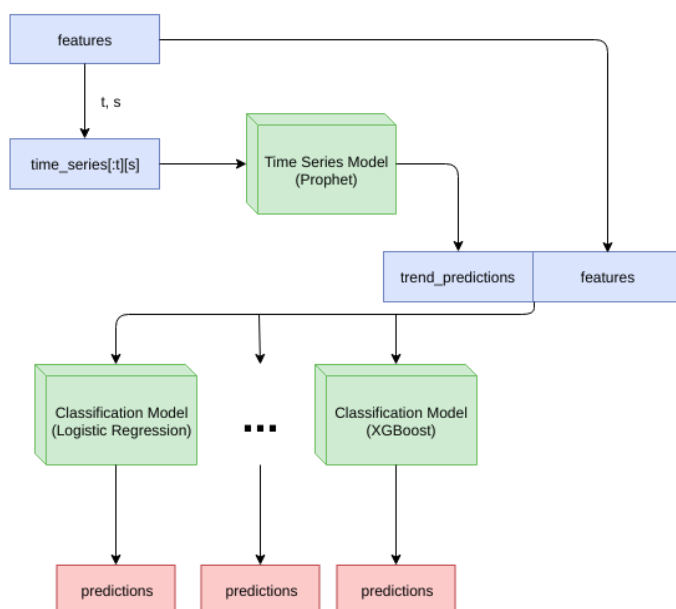
Pipeline and Modeling

In this section, we try to predict gun homicide increases for a state in each month. We have data on gun violence incidents for 2014-2017. This means we'll be training on 2014-2016 and predicting for each of the 12 months in 2017. From the visualization notebook, we saw that there is seasonality and trends in gun homicide rates. In order to use this in a model along with other features, I decided to use a time series analysis algorithm, Facebook's Prophet, to make numeric predictions on gun violence. Here's a visualization of the performance of the time series analysis (more information can be found in the time-series notebook):



I use the output of the time series analysis along with the other features as input into classification models to make predictions on the label. This process is repeated for each state separately. The reason for this decision is that each state has different gun violence patterns and feature effects. Here's a simple diagram of the pipeline:

Making predictions on a row from features
corresponding to **time t** and **state s**



I excluded states with less than 2 average monthly gun homicides. This was because gun homicides is not as significant of an issue in those states, and the label has high fluctuation.

Results and Challenges

I found that XGBoost had the best performance among the models trained. When training on years 2014-2015 and making predictions for 2016, the accuracy was 69%. Training on years 2014-2016 and making predictions for 2017, the accuracy was 75.5%.

Adding just one more year of training increased accuracy by 6.5%. It seems that the model can still learn, but suffers from lack of data.

Furthermore, the features we are using are annual features. This means we won't be able to capture the monthly changes for each feature, and the most useful feature was the prediction for the number of gun homicides from the time series analysis model.

III. Annual Modeling

Pipeline and Modeling

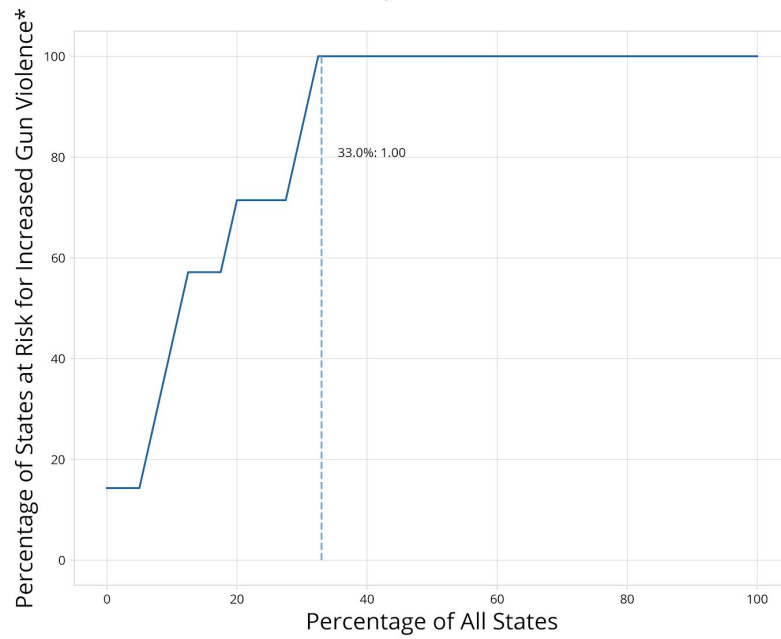
Here, we predict gun violence increases for each state at the annual level. We have gun violence information from 2000-2017, and we will be able to make better use of our annual features here. We have insufficient data points (16 data points from 2000-2016 for training) to be using time series analysis here. We also won't be training a separate model for each state because there are too many parameters to learn from the limited number of data points. Furthermore, we only used the features that were available from 2000 to 2017. This includes features on gun violence rates, gun control laws, and party support for presidential candidates.

Just as with the monthly modeling, I excluded states with less than 20 average annual gun homicides.

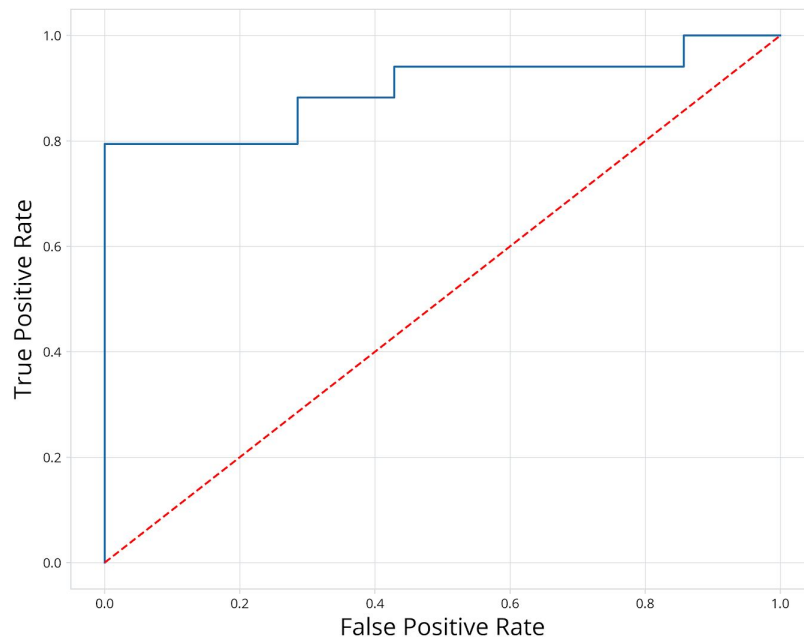
Results and Challenges

Once again, the best performing model was XGBoost. We ended up with 85.4% accuracy with 42.9% recall rate. By tweaking the sensitivity, we got 80.5% accuracy and 86% recall. All of the states that will see a 20% increase in gun homicides were within the top 33% states classified most at risk model. Here are the cumulative accuracy profile and the ROC curve:

Cumulative Accuracy Profile for 2017 Predictions



XGBoost ROC Curve for 2017 Predictions



The total number of laws and laws that restrict access to guns for children had high feature importance in the model. More details can be found in the [annual_modeling notebook](#).

Due to the limited availability of some of the features, we were not able to use all of them. This was unfortunate because for many of these features, we saw correlation with gun violence during the visualization part.

Conclusion

Throughout this project, we were able to find significant correlations between gun violence and gun control provisions. Further exploration of these correlations may lead to a solution to America's growing gun problem. It will be up to lawmakers to decide whether to use data as evidence when deciding on bills.

Our predictive modeling could be significantly improved by more data. Many of the features that we would have liked to use in our models were too limited in terms of available time frames. However, as the quality and quantity of relevant data increases, statistical analysis may become more relevant to policy making.