# Race as an Indicator in Fatal Police Shootings in the United States

Anna Harvey

Bellevue University, College of Science and Technology,

anna.f.harvey@gmail.com

**Summary**

The Washington Post has collected information about fatal police shootings from the news, social media, and police reports since 2015[1]. The efforts of such bodies to collect this type of data is necessary to understanding the full extent of police interactions with the public. However, collecting information only about the fatal police shootings and not on non-fatal shooting or other confrontations cannot provide enough information to predict the outcome of police altercations.

**Business Problem**

Police violence in the United States has been a subject of intense discussion over the past several years. The information the public receives about police shootings comes through various media sources that present information with an inherent bias. When discussing what motivates police shootings, going to the raw data for analysis may help mitigate biased interpretation and paint a better overall picture. When media or politicians report on police shootings, the bias of the individual or organization may misrepresent data to the public to make an argument for their position. Extracting the raw data to show the basic facts of each situation outside of political or personal interpretation can allow more appropriate policies to be created and enforced to reduce the number of fatalities from police/citizen interactions.

**Project Proposal**

I will use the US Fatal Police Shootings (2015-current) dataset from the Washington Post to investigate correlations between race, perceived mental illness, age, gender, and whether the person was armed. The goal will be to discover if there are any strong correlations between the person's attributes and being fatally shot by police. I acknowledge the fact that these scenarios are complex, and the dataset leaves out information such as previous criminal charges, and the nature of the confrontation. The data may itself contain bias due to the nature of data collection from news sources and social media. Official fatal police shooting reports are also voluntary and it is probable that this dataset does not include all fatal shootings from the last five years and may include inaccurate or missing data. Without strict governmental oversight to the collection of accurate police reports, the data presented will be sufficient for exploratory analysis.

The primary goal for predictive analysis is to determine whether or not a person's race will predict the outcome of an altercation with police and if other factors contribute to the likelihood of that prediction. Therefore, the focus will be on multi-target classification models.

---

[1] *Fatal Force: Police shootings database*. (2020, January 22).
https://www.washingtonpost.com/graphics/investigations/police-shootings-database/.

**Implementation**

I was able to conduct some exploratory data analysis on the dataset, removing irrelevant details for my problem statement, such as location, name, and manner of death. Initial overview shows a dispersion of altercations across seven race variables, including unknown race and "other" (Fig. 1). My final dataset to use for analysis included Race, Age, Gender, Mental Illness (whether or not the victim was reported has having signs of or a diagnosed mental illness), and Armed (whether or not the victim was reported as being armed).
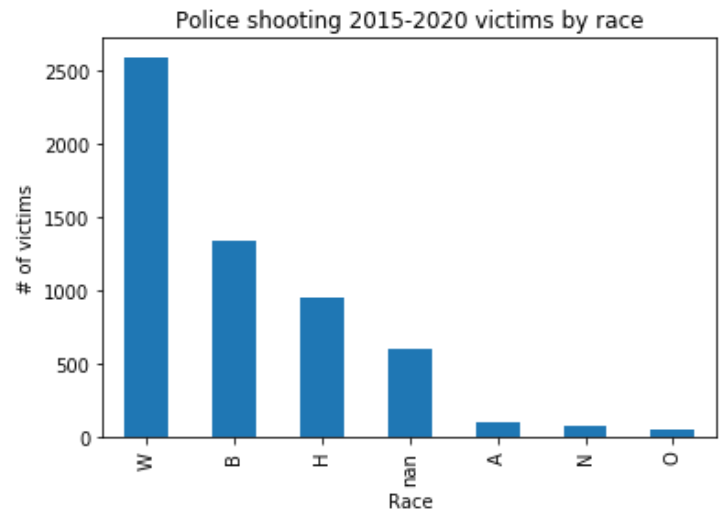


Fig. 1: Distribution of altercations by Race. W = white, B = black, H = Hispanic, nan = Unknown, A = Asian, N = Native American, O = Other

Data was split into test sets and training sets before conducting analysis meant to infer information from the data. I attempted to use One-Hot Encoding to encode the race variables, as they are categorical but not nominal. However, this did not work well for further analysis, as race is my target variable. I returned to the original data set and used label encoding instead to transform the race variables into numeric values (0-6). I also encoded the other categorical variables (Gender, Mental Illness). After this, I was able to establish a Pearson's Correlation matrix to analyze potential relationships (Fig. 2).
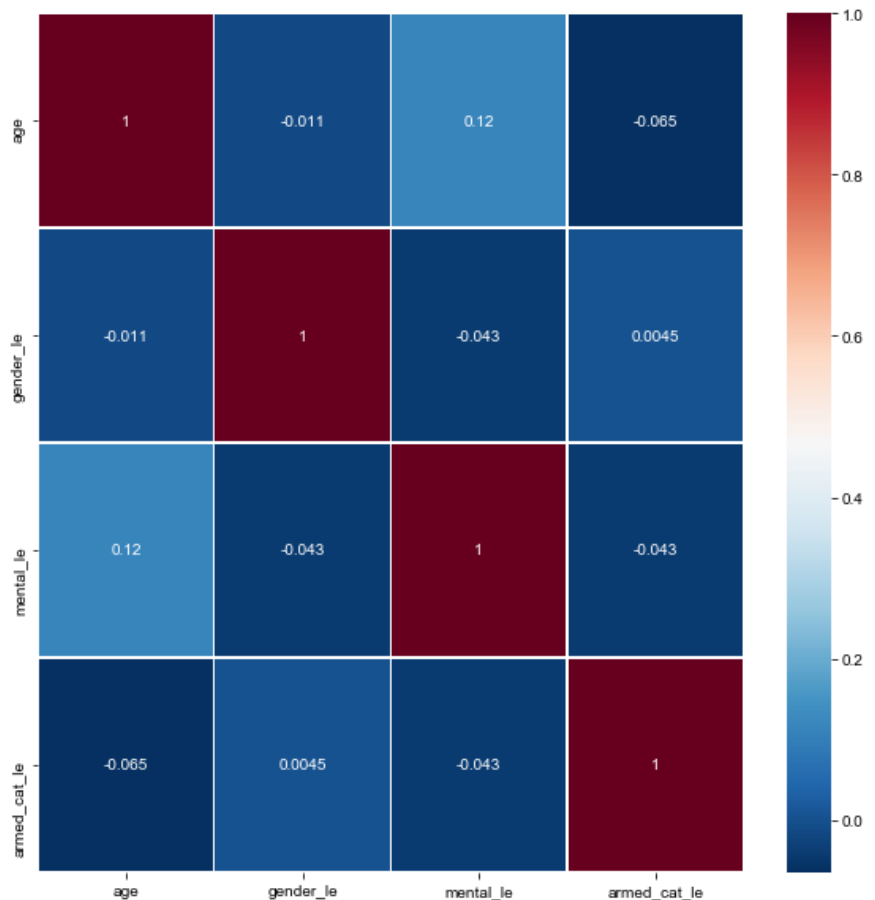


Fig. 2: Pearson's Correlation matrix. The strongest correlation appears to be between Age and Armed.

An attempt was made to use Principle Component Analysis (PCA) to analyze whether any variables were unnecessary to the analysis. This method worked well when the data was One-Hot encoded as the variables of Gender, Mental, and Armed were all binary variables. However, with a dataset with only four variables, it was not necessary.

Multi-target classification models of Naïve Bayes, K-means, and Random Forest were used to construct classifiers. K-means was not an appropriate method to use for this dataset as it does not sort the data based on the intended labels of my target variable. For further study, I would require more data and more specific types of classification models to handle categorical data.

**Results**

Due to the nature of the dataset, I was unable to reach any conclusive results. What I have learned from this project is that to apply neural networks and machine learning models to data, there needs to be a clear target variable and enough data. Understanding the metrics required to evaluate the efficacy of the model is crucial.

**Conclusion**

To accurately predict whether an altercation with police will end in a fatal shooting based on someone's race, a dataset regarding non-fatal police shootings would need to be integrated into analysis as well. With only a dataset of fatal police shootings, we already know the outcome of the altercations and there are many other variables to be considered during police confrontations. To thoroughly investigate police violence from an unbiased perspective, extensive datasets would need to be created and maintained with details regarding all police confrontations, regardless of outcome. It is a large task to undertake and would require more direct and accurate police reporting, as well as a national requirement for such reporting to take place. The Washington Post dataset is a good resource for exploratory data analysis but more thorough data is required for predictive data mining.