Policing the Police

# Our Dataset

**Merged from the largest databases of police violence in America:**
Mapping Police Violence, Deadspin, The Washington Post

**Consists of all recorded police killings from 2013 - 2020**

**Individual-Specific Variables**
- Age
- Gender
- Race

**Location-Specific Variables**
- State
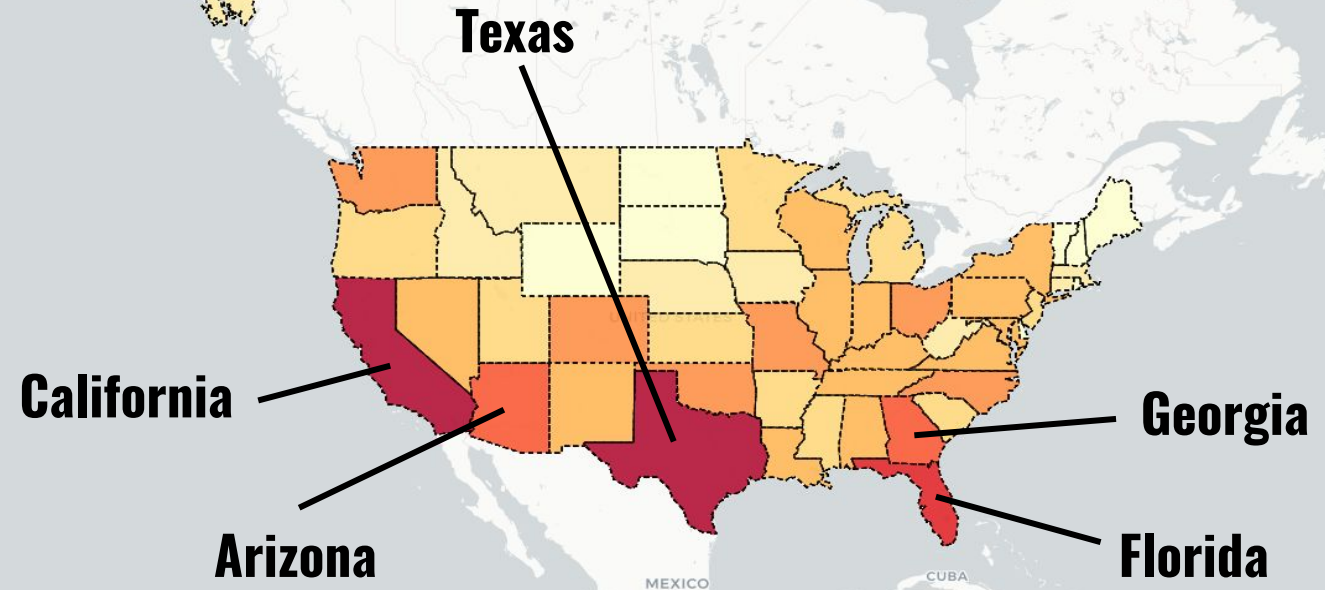- City
- County
- Geography type: Rural/ Urban/ Suburban

**Circumstance-Based Variables**
- Weapon on victim
- Whether the victim was perceived to be a threat
- Whether the victim suffered from mental illness
- Cause of death
- Description of death
- Whether the police killing was justified under law
- Police agency responsible

Rate of Police Killings (Per Million)

Top 5 States

Shootings per million people

0  1  3  4  6  7  9  10

Alaska

Nevada

Arizona

New Mexico

Oklahoma

Leaflet | © OpenStreetMap contributors © CartoDB, CartoDB attributions

Rate of Police Killings (Per Million)

Bottom 5 States

Shootings per million people
0  1  3  4  6  7  9  10

New York
Massachusetts
Rhode Island
Connecticut
New Jersey

Leaflet | © OpenStreetMap contributors © CartoDB, CartoDB attributions

# How Have Police Killings Changed Over Time?



Police Killings of Top 5 States over the years
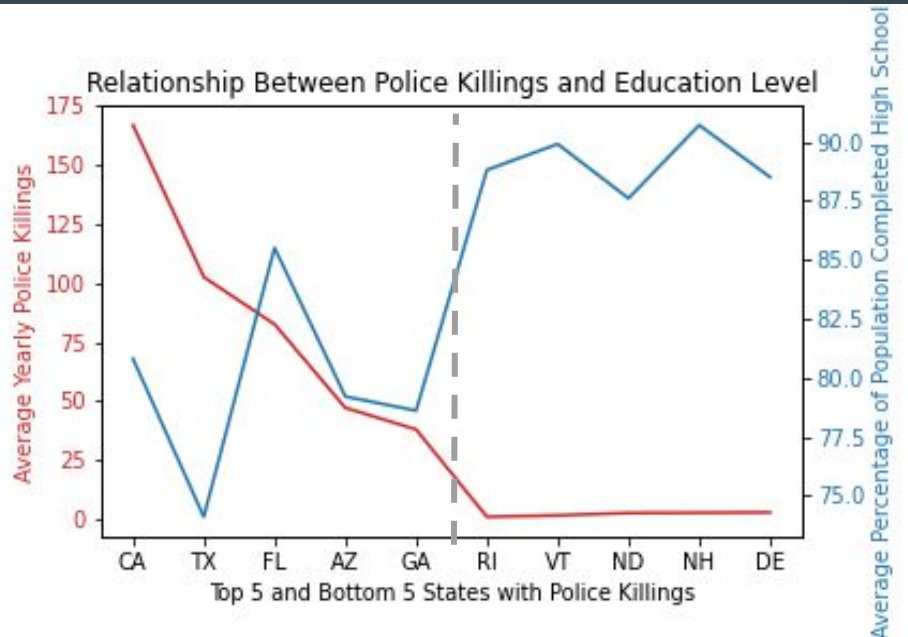


Police Killings of Bottom 5 States over the years

Small fluctuations over the years: Police violence has not improved in these Top 5 states!

Bottom 5 states have stayed relatively low over the years
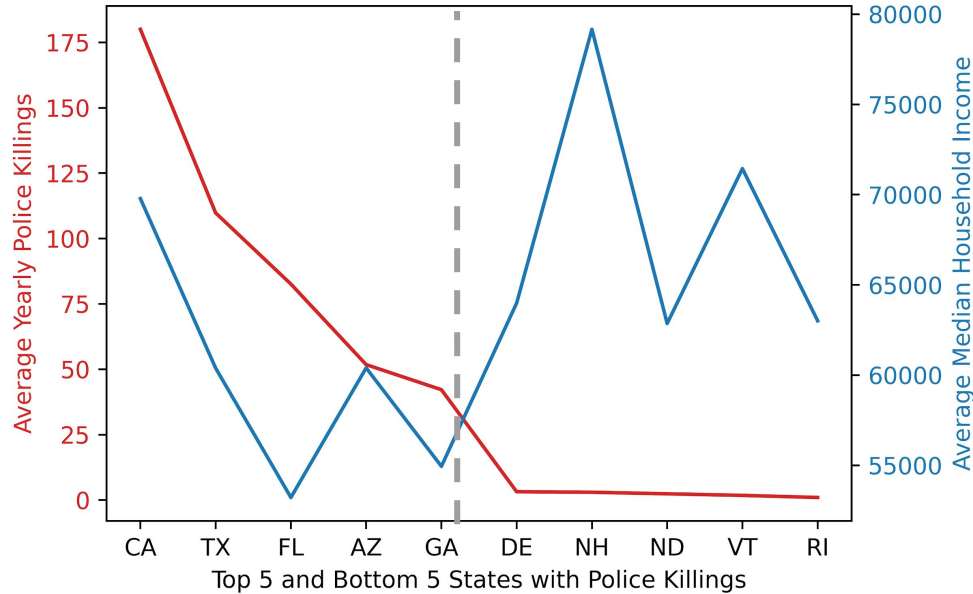
# Why Do States Differ In Police Killings?



Relationship Between Police Killings and Education Level

**Negative Correlation:**
Higher education level = Lower police killings

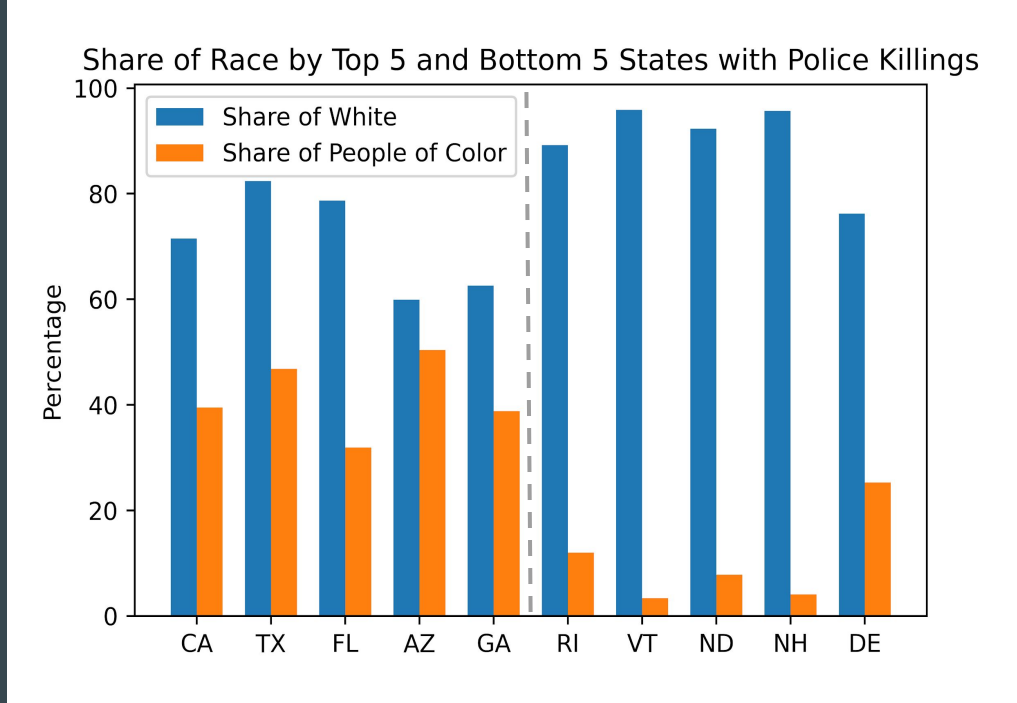## Different Education Levels

# Why Do States Differ In Police Killings?



Relationship Between Police Killings and Median Household Income

Average Yearly Police Killings

Average Median Household Income

Top 5 and Bottom 5 States with Police Killings

CA TX FL AZ GA DE NH ND VT RI

**Negative Correlation:**
Higher income level = Lower police killings

## Different Income Levels

# Why Do States Differ In Police Killings?



Share of Race by Top 5 and Bottom 5 States with Police Killings

Legend:
- Share of White
- Share of People of Color

X-axis states: CA, TX, FL, AZ, GA, RI, VT, ND, NH, DE
Y-axis: Percentage (0 to 100)

**Positive Correlation:**
Higher share of people of color
= Higher police killings

## Different Population Racial Makeup

# Why Do States Differ In Police Killings?

Create 3 new state-level features:
1) ***state_education_level:*** Proportion of population completed high school in each state
2) ***state_log_income:*** Log of median household income in each state
3) **state_white_share:** Share of Whites in each state

Creating New Features Using Correlations Observed

# Circumstances Surrounding Police Killings



Circumstance 1: Police killings inside residences

# Circumstances Surrounding Police Killing



Circumstance 2: Police killings involving vehicles, where victim is fleeing

# How Were Victims Killed?

Group into 3 most common causes of death:
- By gunshot
- By taser
- By physical violence (Beaten/ Physical Restraint/ Asphyxiation/ etc)

Create 3 new binary features:
1) *killed_by_gunshot*
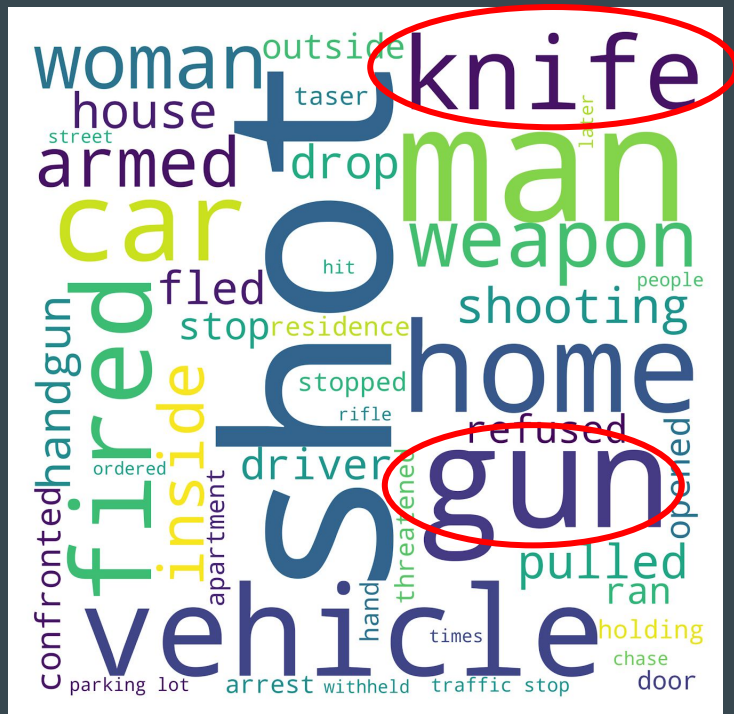2) *killed_by_taser*
3) *killed_by_physical_violence*

## Analysis of *Cause of Death* Column

| | |
|---|---|
| Gunshot | 6982 |
| shot | 1437 |
| Gunshot, Taser | 237 |
| Taser | 221 |
| shot and Tasered | 88 |
| Beaten | 29 |
| Vehicle | 27 |
| Physical Restraint | 23 |
| Tasered | 13 |
| Physical restraint | 9 |
| Asphyxiated | 8 |
| Gunshot, Police Dog | 5 |
| Other | 5 |
| Pepper Spray | 4 |
| Taser, Physical Restraint | 2 |
| Gunshot, Pepper Spray | 2 |
| Taser, Pepper spray, beaten | 1 |
| Baton, Pepper Spray, Physical Restraint | 1 |
| Bean bag | 1 |
| Beaten/Bludgeoned with instrument | 1 |
| Bomb | 1 |
| Chemical agent/Pepper spray | 1 |
| Gunshot, Beanbag Gun | 1 |
| Gunshot, Bean Bag Gun | 1 |
| Gunshot, Stabbed | 1 |
| Taser, Pepper Spray, Beaten | 1 |

# Data Transformation to Create Meaningful Features

# Were Victims Armed?

**Analysis of _Description of Death_ Column**



**Analysis of _Weapon On Victim_ Column**

| | |
|---|---|
| Allegedly Armed | 5383 |
| Unarmed/Did Not Have an Actual Weapon | 1082 |
| gun | 886 |
| Unclear | 612 |
| Vehicle | 505 |
| knife | 229 |
| unarmed | 103 |
| toy weapon | 45 |
| undetermined | 40 |
| vehicle | 35 |
| unknow | 20 |
| sword | 9 |
| machet | 9 |
| baseba | 7 |
| metal | 7 |
| gun an | 6 |
| hammer | 5 |
| hatchet | 5 |
| ax | 4 |

We already created a new feature named v*ehicle_involved* earlier

Create 2 new binary features:
1) *armed_with_gun*
2) *armed_with_knife*

## Data Transformation to Create Meaningful Features

# Target Encoding: For Categorical Features With Many Levels

Categorical features encoded:
1) *state:* **Consists of 50 unique states**
2) *police_agency:* **Consists of 2866 unique police agencies**

**What is Target Encoding:**
Encode each level with the mean of the target variable for that level

**Advantage over One-Hot Encoding:**
Does not add to dataset dimensionality

Encoding Categorical Variables

# Predicting the Race of Individuals Killed by Police

**Dependent Variable (Y):** *Black*
Shows whether an individual is Black (0 or 1)

**Hypothesis**
Given that an individual was killed by the police, there is sufficient difference in the manner they were killed to differentiate whether they are Black
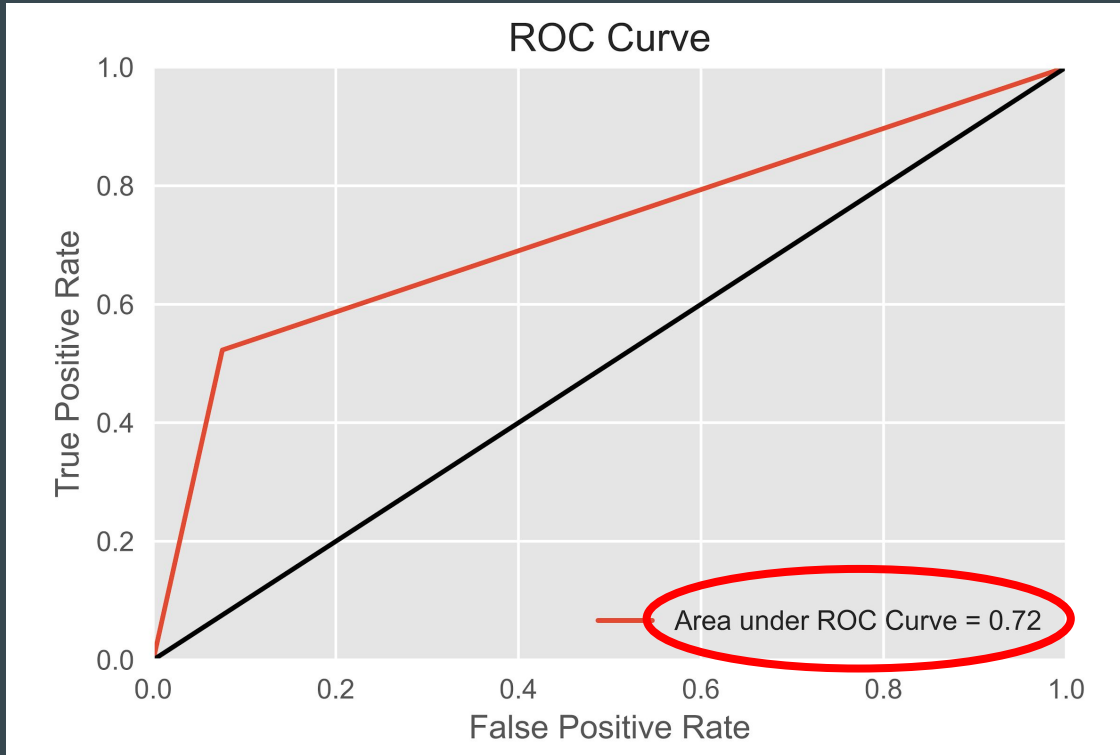

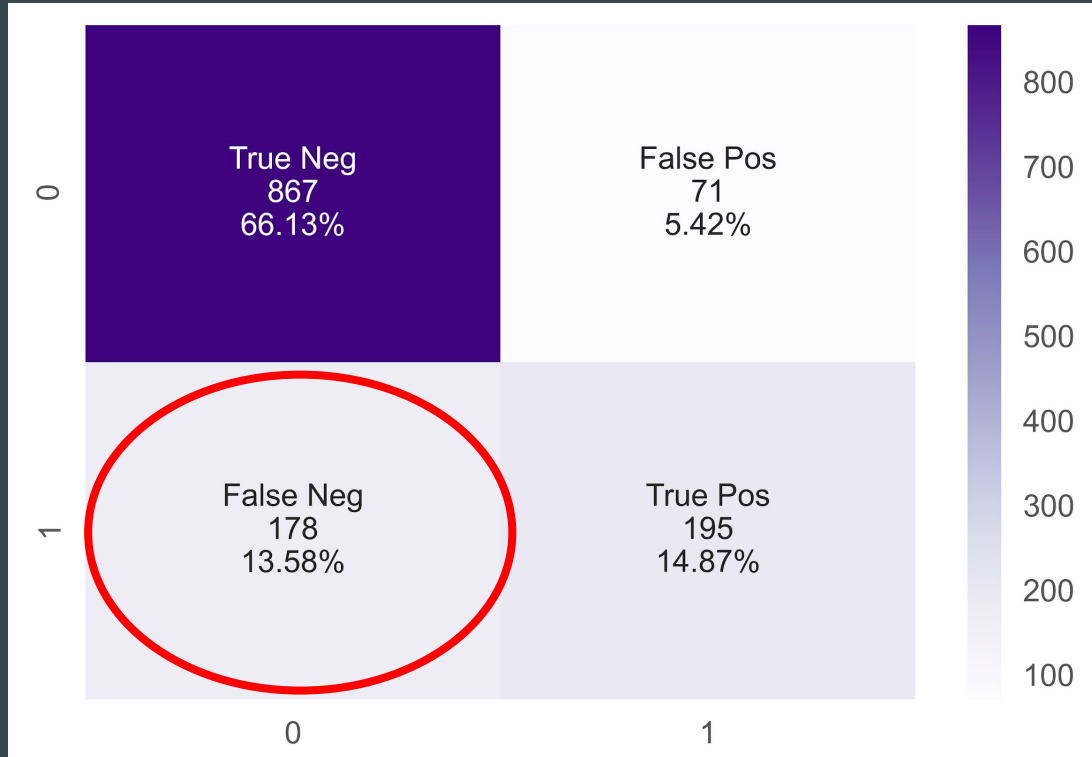
Correlation of Y Variable, "black", With X Variables
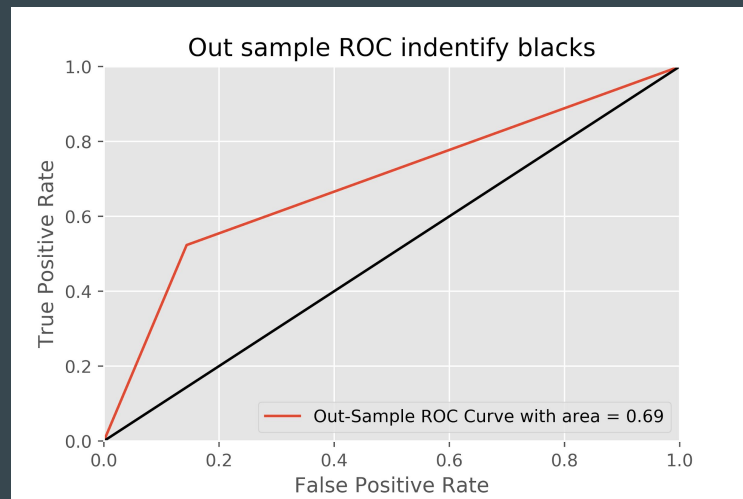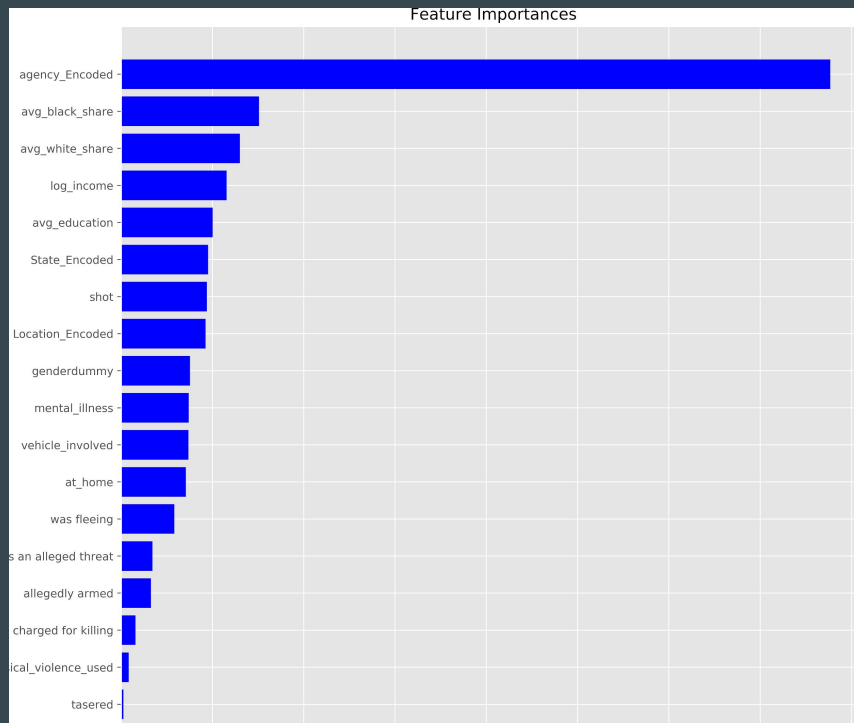
# Predicting the Race of Individuals Killed by Police

**Advantages of Random Forest Model:**

1) Does not assume linearity between Y and X variables

2) Able to handle missing values:
   Column ***justified*** has over 4000 rows with missing values
   -> We can include it in our Random Forest model but could not in our
   Logistic Regression model

## Model 2: Random Forest

# Predicting the Race of Individuals Killed by Police
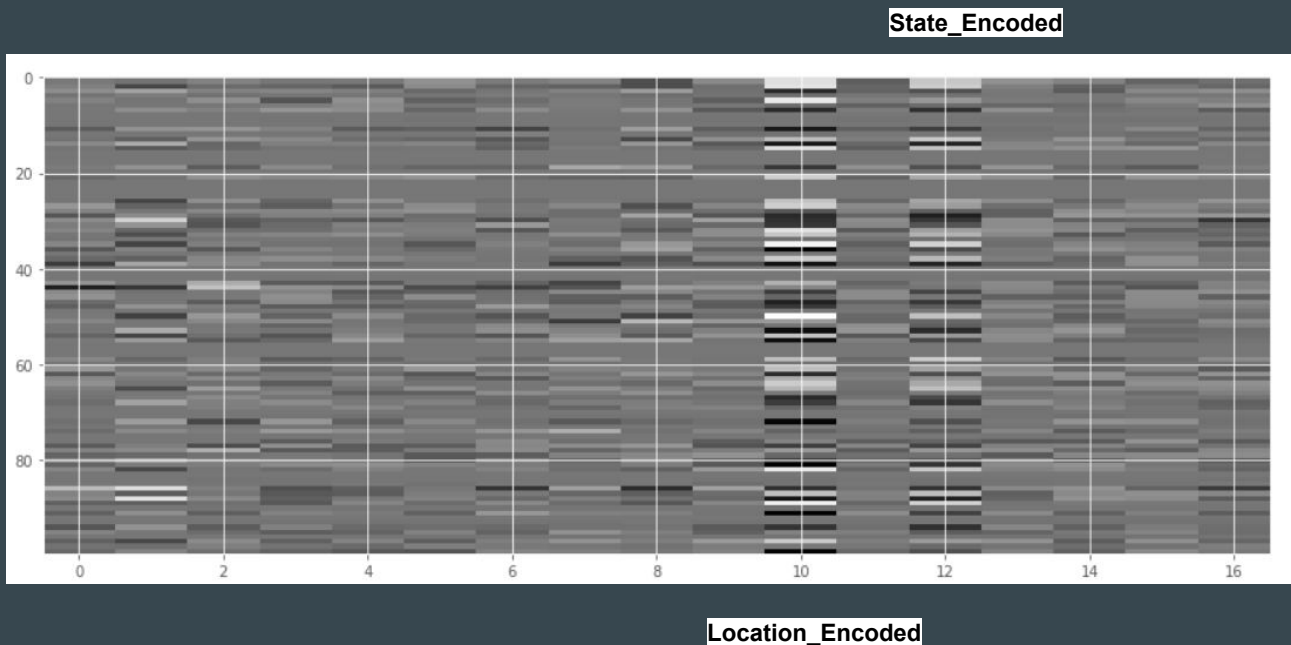


**Accuracy Score: 0.74**

CFM: [825, 126],
[178, 189]

## Model 2: Random Forest

# Predicting the Race of Individuals Killed by Police



**Accuracy Score: 0.81**

CFM: [868,  60],
     [260, 130]

## Model 3: Neural Network