



Final Project

Laura Godleski
Lindsay MacDonald
Natalie VanDyke



Overview

WHO IS OUR TEAM?



Lindsay
MacDonald



Natalie
VanDyke



Laura
Godleski

WHAT IS THE TOPIC?

*Influence of
demographics on
the homicide solve
rate in California.*

WHY THIS TOPIC?



Data Sources

U.S. Homicide Reports, 1980-2014

<https://www.kaggle.com/jyzaguirre/us-homicide-reports>

US Census Demographic Data

https://www.kaggle.com/muonneutrino/us-census-demographic-data/data?select=acs2015_county_data.csv

Questions to Answer



- Is there a correlation between the personal attributes of the victim and the solve rate?
- Is there a correlation between county demographics where the murder occurred and the solve rate?
- Can we develop a machine learning model that predicts whether or not a crime would be solved given hypothetical sets of circumstances?

Data Exploration: Preprocessing

Kaggle Dataset #1 (US Homicide Reports, 1980-2014)

- Filtered 'State' column for California
- Renamed 'City' column as 'County' in order to merge with the second dataset
- Dropped irrelevant columns
- Checked for null values (0)

Kaggle Dataset #2 (US Census Demographic Data)

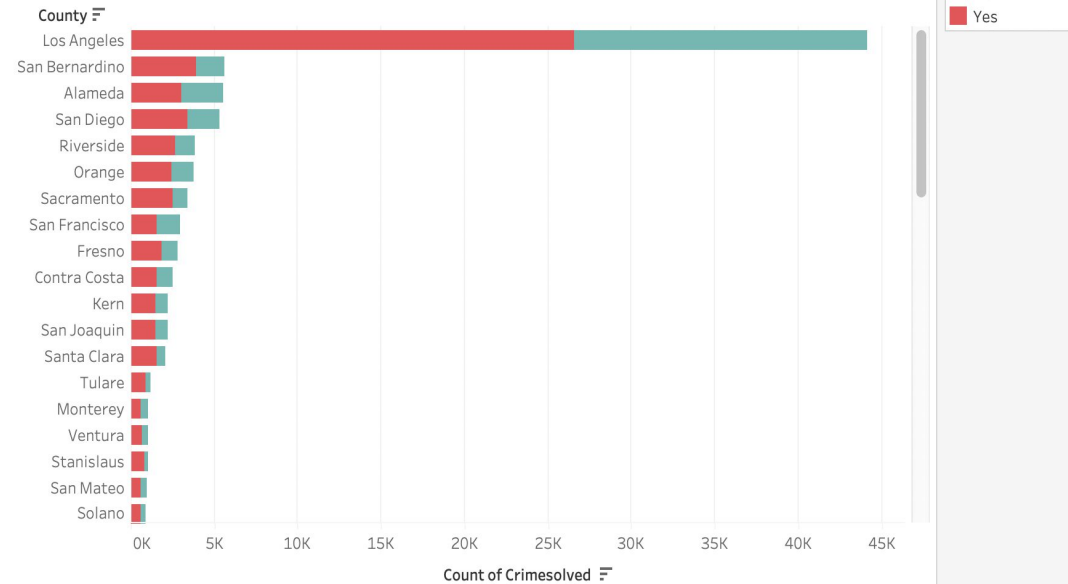
- Filtered 'State' column for California
- Dropped irrelevant columns
- Checked for null values (0)

Datasets were imported into Postgres and joined; merged dataset was then uploaded to AWS S3 bucket

Data Exploration: Initial Findings

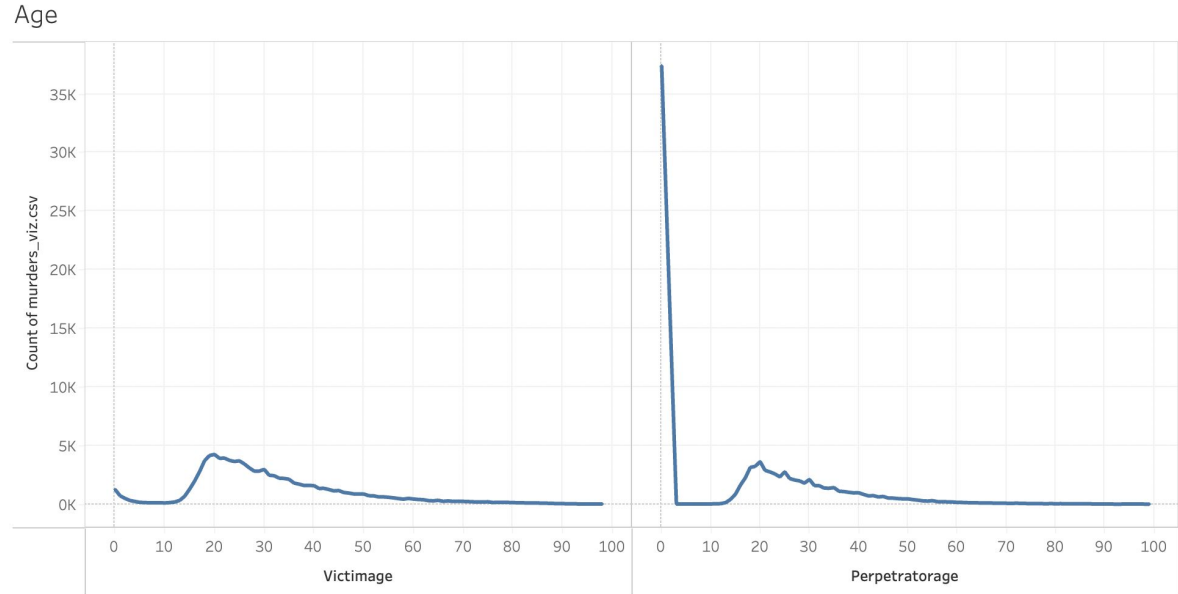
Due to the population disparity among counties in CA, a vast majority of the data was for Los Angeles county

Solve rates by County



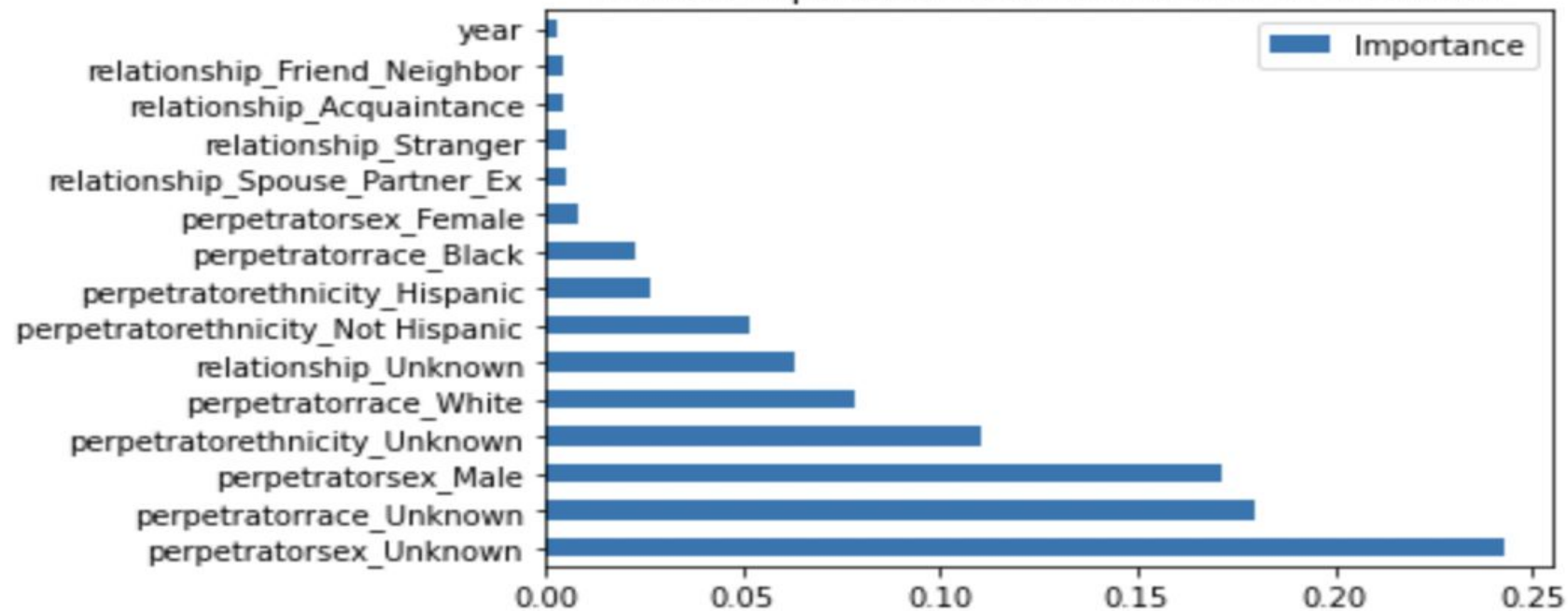
Data Exploration: Initial Findings

There is an unusual amount of records capturing an age of 0 for perpetrators which suggests these individuals were largely unidentified and therefore their age was unknown



Analysis

15 Most Important Features in Outcome Prediction



Analysis

