# Re-evaluating Recidivism

Benjamin Sperling

# Recidivism: What is it and what are the stakes?

- Recidivism is the likeliness of a person convicted of a crime to reoffend
- This is a rather large umbrella term, as convicted criminals vary in the magnitude of their crimes
- Prisons use COMPAS, an algorithm that looks at many variables of a convicted person and determines the likelihood of reoffending
- Predicting recidivism can potentially save lives by keeping dangerous criminals in jail
- However, the COMPAS algorithm has been shown to be inherently biased against minority groups
- COMPAS should be re-evaluated to examine its effectiveness and efficiency
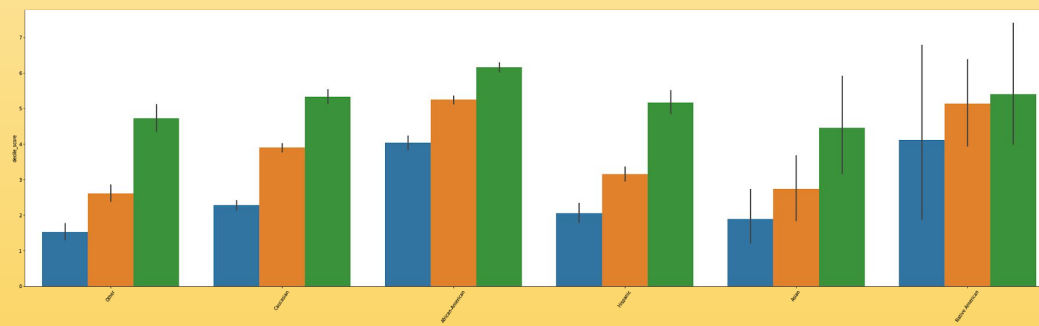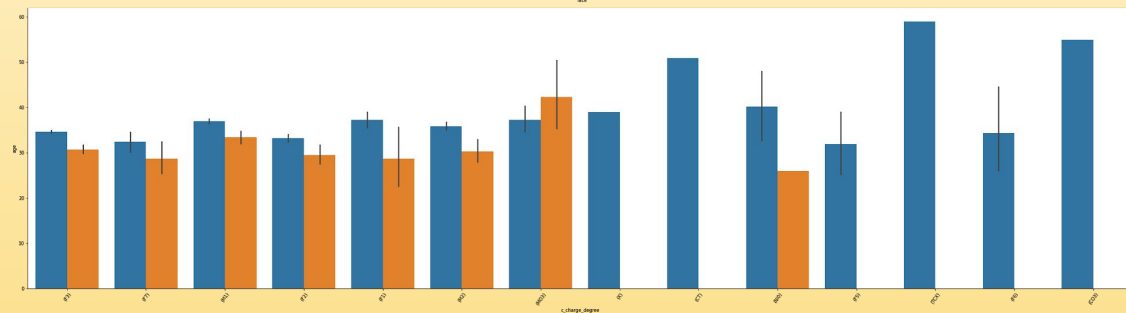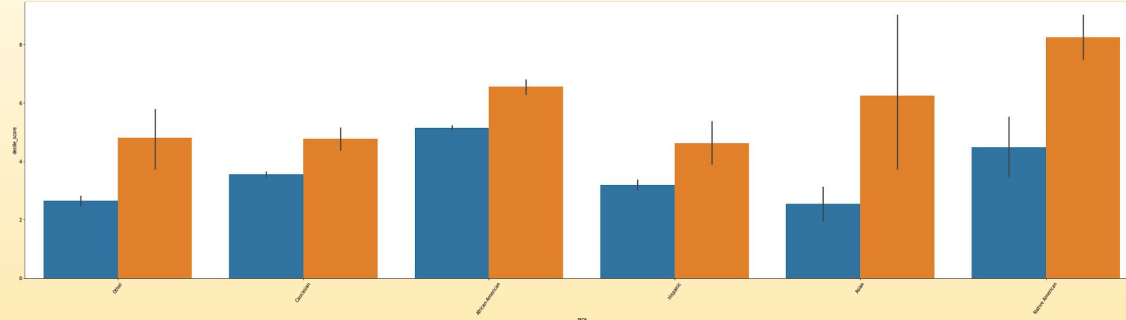
# Hypotheses/Questions

- Hypotheses: Sex, Age, Time Spent In Jail, and Prior Offenses are the best ways to determine risk of recidivism based on the given Violent-Parsed dataset
- Important: Model should be trained on subset of data for violent recidivism
- Avoiding re-incarceration due to violence should be the number one priority
- Including non-violent re-offenses would skew data
- How does age factor into recidivism?
- Does the amount of time spent in jail make the occurrence of repeat offenders more or less likely?
- What other data would be useful in creating a model to calculate risk of recidivism?
- How do we approach race as a data point?

# EDA

- Created a new dataset with just the independent variables and dependent variable
- Independent Variables
  - Sex, Age, Time Spent In Jail, and Prior Offenses
- Dependent Variable
  - Violent Decile Score
- Graphed relationships to identify trends between variables
- Changed certain STR variables to FLOAT to allow for better modeling

# EDA (cont.)

1. Violent recidivism vs. non-violent recidivism by race
2. Violent recidivism vs. non-violent recidivism by charge degree
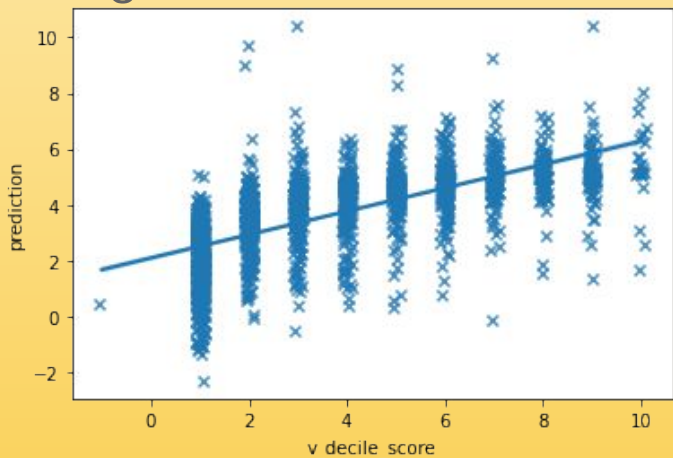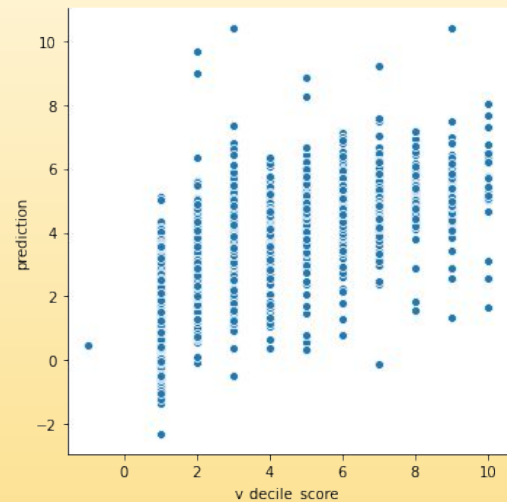3. Decile score and race by age category

# Model

- Simple Linear Regression Model
- Established Cross-Variance
- Key Variables to base the model off
  - Sex
  - Age
  - Time In Jail (first offense)
  - Juvenile/Prior Offenses
- Nuanced data to begin with, specifically in regards to race (more in conclusion)

# Results

- Prediction model is somewhat accurate
  - Low amount of variables = less accurate
- Line of best fit does trend in the right direction
- Notable outliers; model is not perfect
- Prediction score goes into negatives
  - Less harmful?

# Key Takeaways

1. A model to predict recidivism should be trained on a very large set of variables
2. These variables are incredibly nuanced
3. These variables are inherently biased due to the need of human correction for specific variables
4. After a COMPAS score is given, it should be used as a guide, not a mandate, when evaluating prisoners
5. Humans are not datasets

# Conclusion

- Data is nuanced, so many other factors go into why certain people have certain scores
- To try to create numerical coefficients to offset these racial and societal advantages/disadvantages is asking to put those categories into ordered lists
- The people doing that categorization will be inherently biased, and based on the current U.S. prison system, inherently biased against minorities and low-income backgrounds
- There should be two separate models used to give two separate scores: one for risk of violent recidivism, and one for risk of non-violent
- More categories should be looked at to train the model