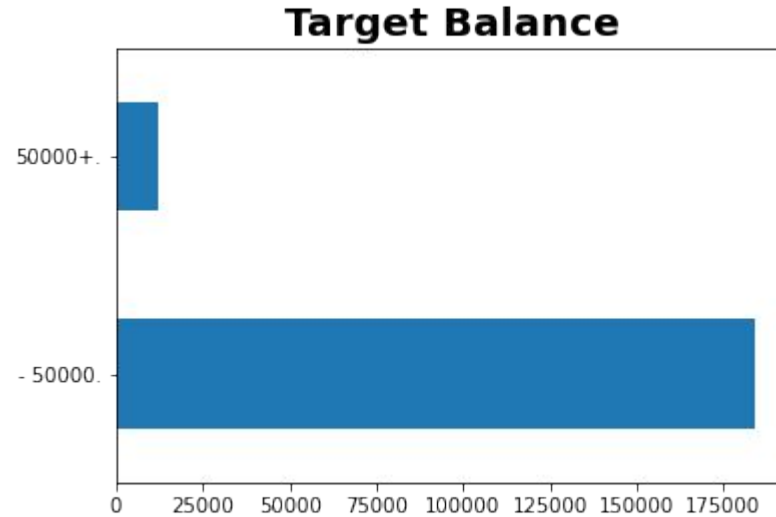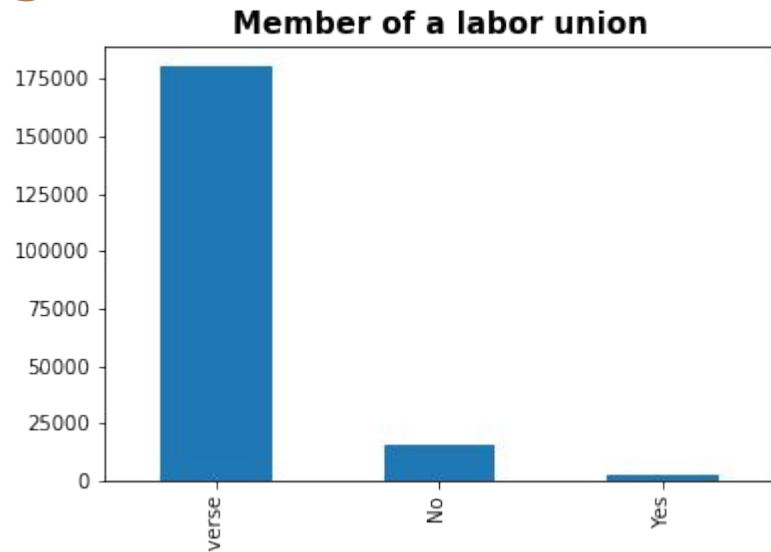# Census Income

Sibongile Toure

# EDA

- Unbalanced Dataset
- Features
  - 33 Nominal
    - Marital Status
    - Class of Worker
  - 7 Continuous
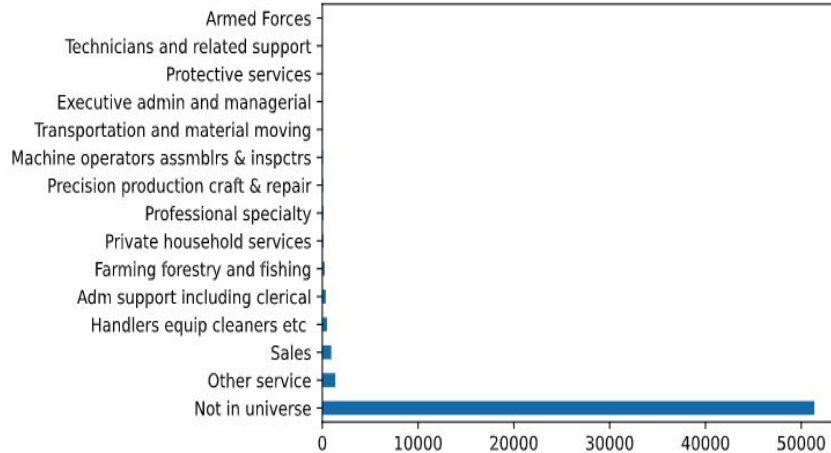    - Age
    - Number of Hours Worked

# Feature Engineering

- Dropped features with ~50% NaN values
- Excluded data for people 18 and Under
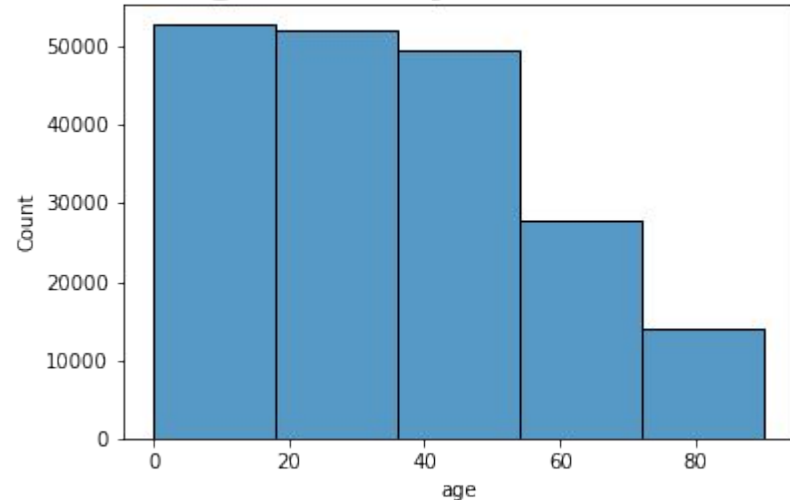


**Member of a labor union**
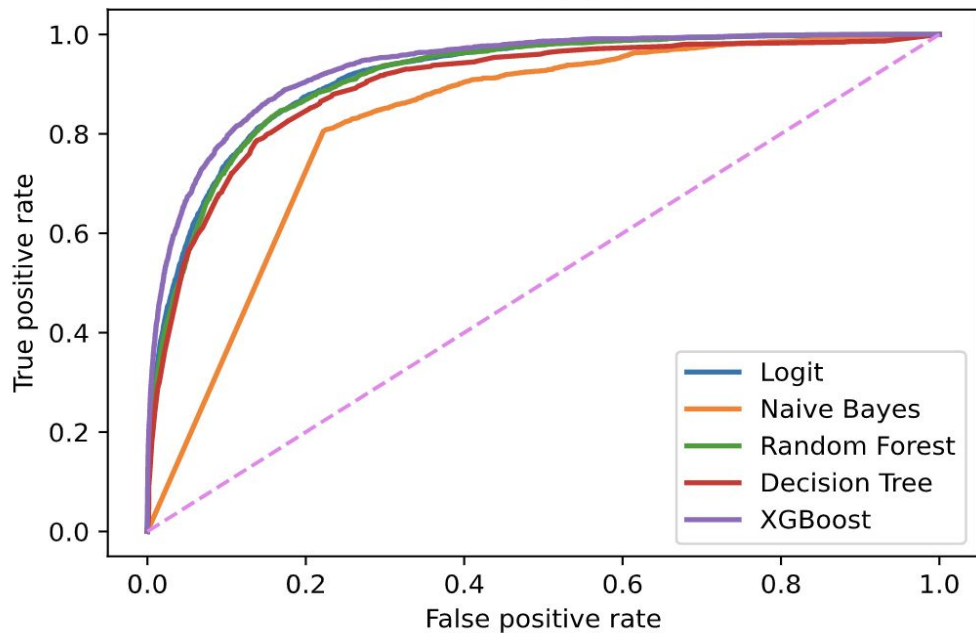
# Feature Engineering



Occupation for people 18 and under



Ages in 18 yr increments

# Modeling

## ROC Scores



## LR + ROS Confusion Matrix

# XGBoost



**Feature Importance**

weeks worked in year
sex_ Male
sex_ Female
detailed occupation recode
dividends from stocks
capital gains
tax filer stat_ Nonfiler
education_ Masters degree(MA MS MEng MEd MSW MBA)
major occupation code_ Other service
detailed household summary in household_ Householder
major occupation code_ Professional specialty
major industry code_Unknown
class of worker_Unknown
education_ High school graduate
detailed household summary in household_ Child 18 or older

**ROC Score**
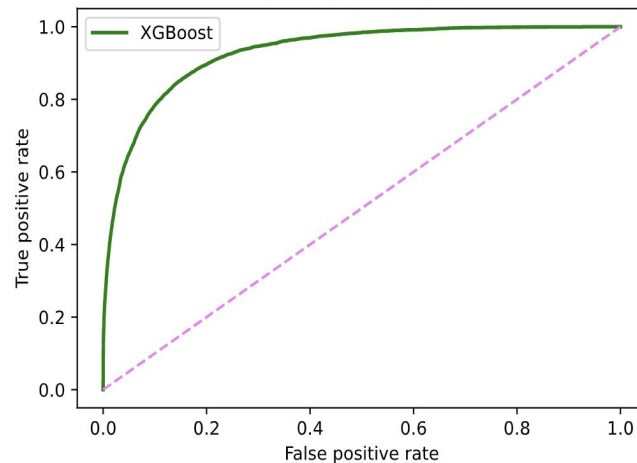
XGBoost

# Takeaways

- By segmenting out people under 18 we are able to get a clearer picture of the workforce
- Certain features such as how many weeks a person worked in a year and their sex has an impact on whether they make above or below $50000

# Future Work

- Include more data from other sources
  - E.g. Neighborhood data can provide insights
- Try additional models