# Applying Supervised Machine Learning Models
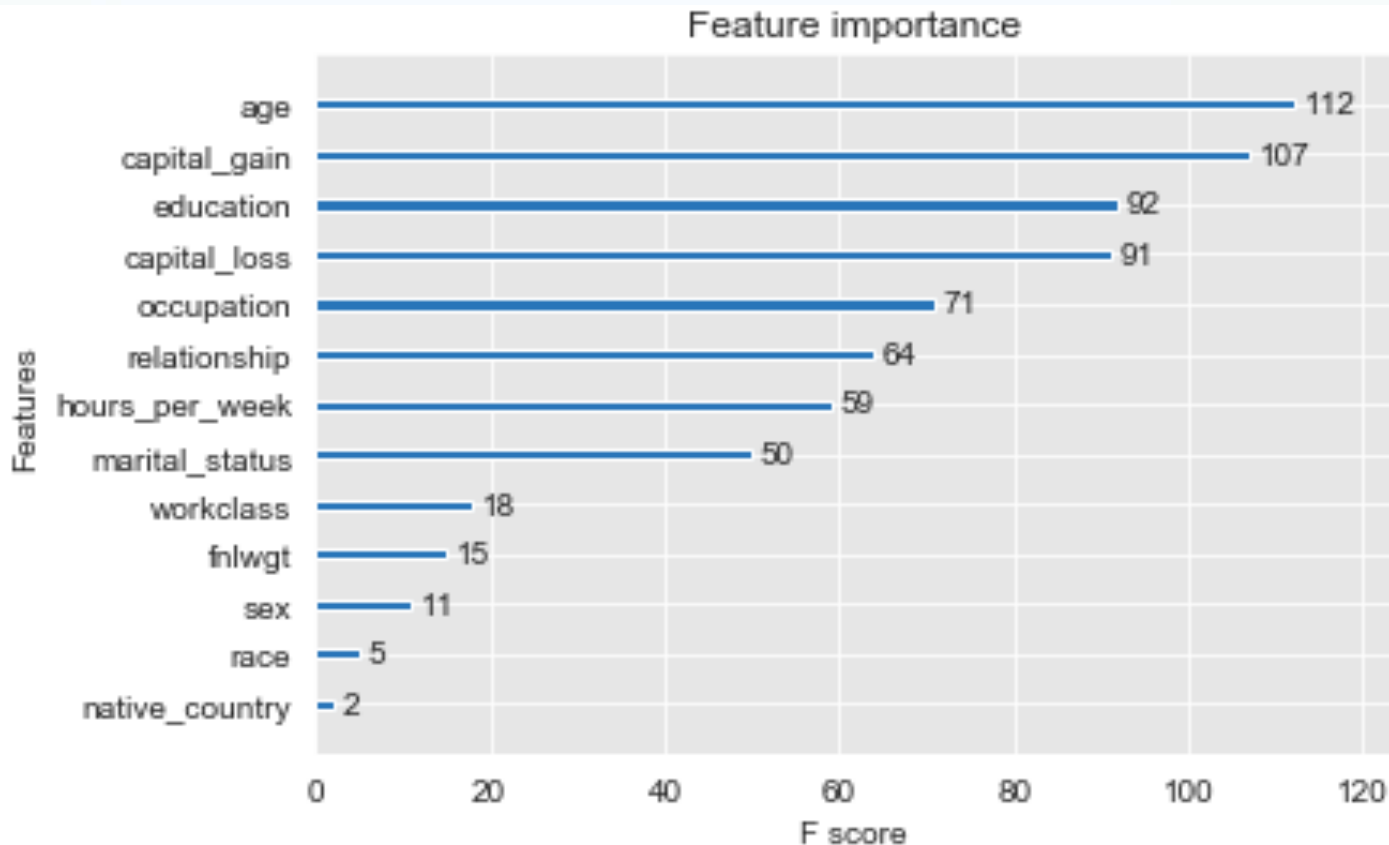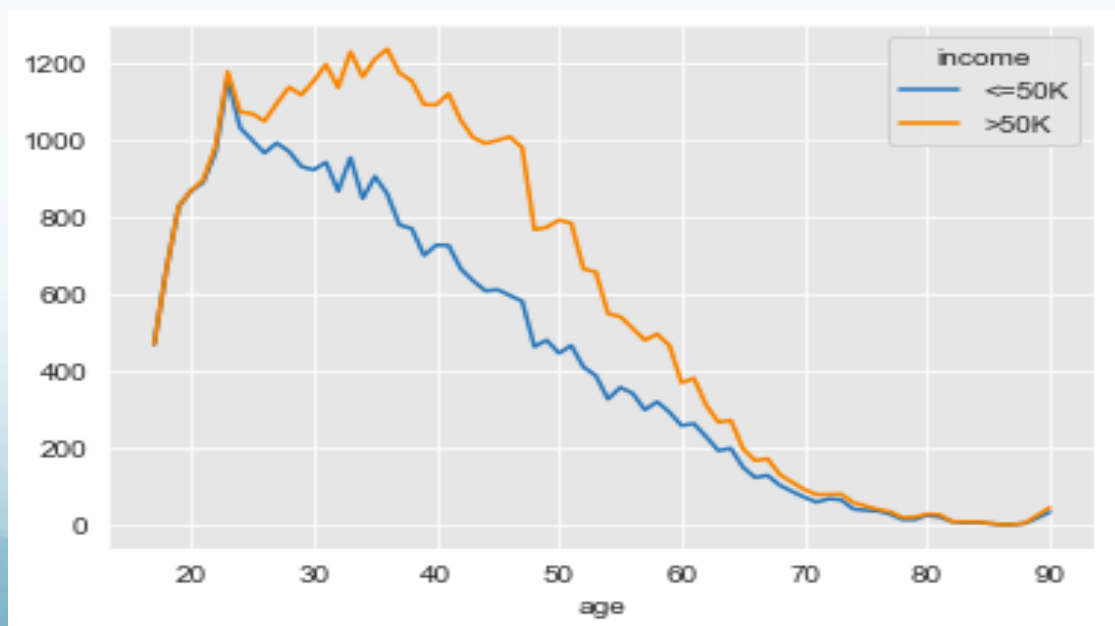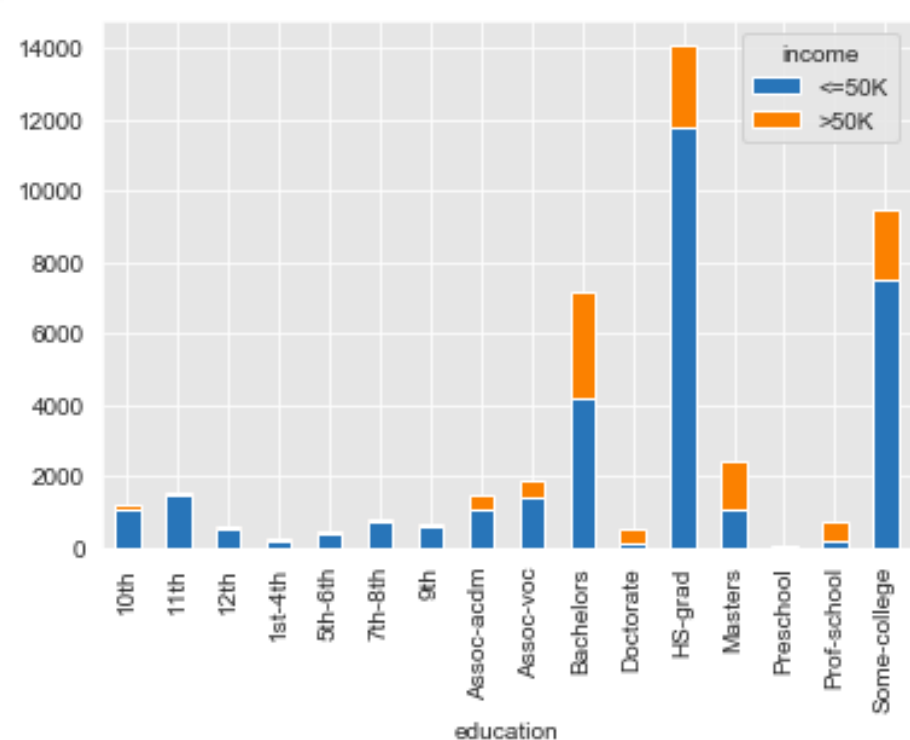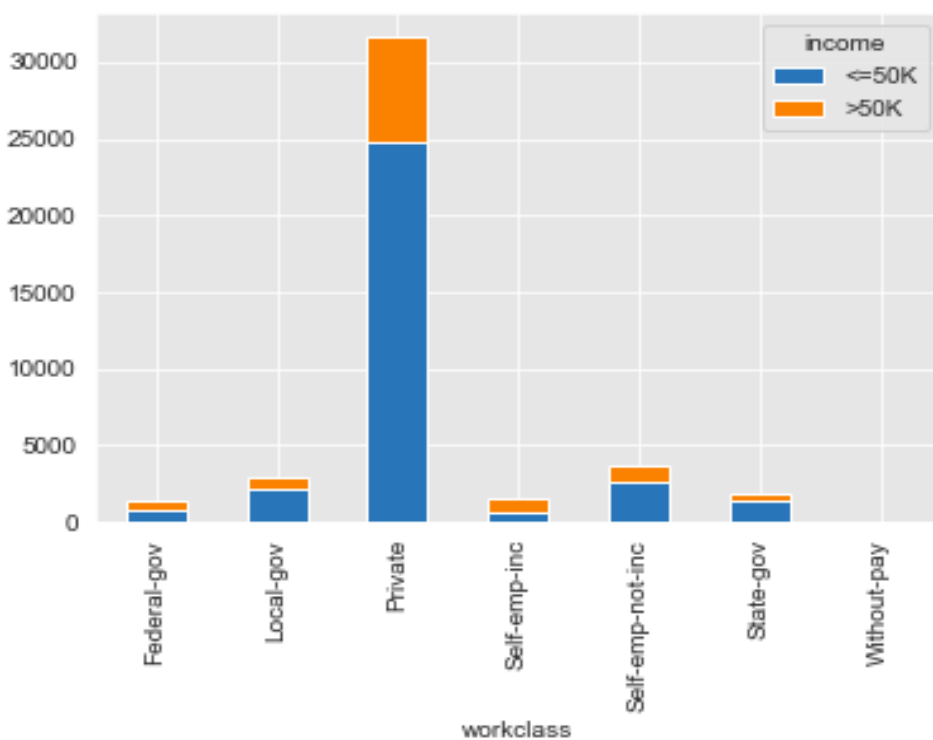
By Natalia Quintero

# Data Set

- Census Income data set from the 1994 census data

- Prediction task is to determine whether a person makes over 50K a year

- Machine Learning problem: classification

- The data set contains 14 attributes and 48,842 instances

# OSEMN Methodology

- **Obtain:** gather information, obtain the data [http://archive.ics.uci.edu/ml/datasets/Census+Income](http://archive.ics.uci.edu/ml/datasets/Census+Income))

- **Scrub:** remove data that is not needed, reduce noise

- **Explore:** set up the data, make sure the dataset meets what is necessary for the type of model to apply later on

- **Classification Models:** logistic regression and XGBoost

- **Interpret**: compare models and evaluate the results

# Important Attributes

# Model Results

- Logistic Regression Model: accuracy of 79%

- XGBoost Model: accuracy of 86%

- Better model for predictions: XGBoost

# Interpretation of Results

- We can predict and classify income with the given data.

- The private sector is were most people work with higher salaries.

- Higher education means higher salaries.

- The peak income for those earning more than 50K is in between early 20's and late 30's.

Thanks for your time