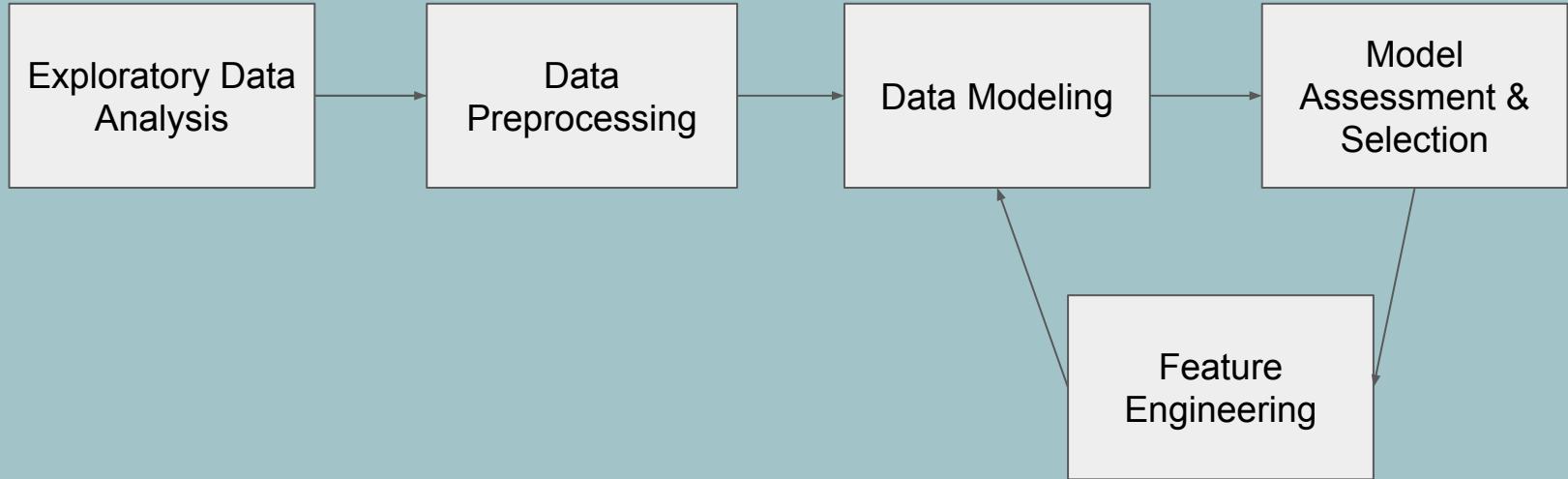


Income Classification Report

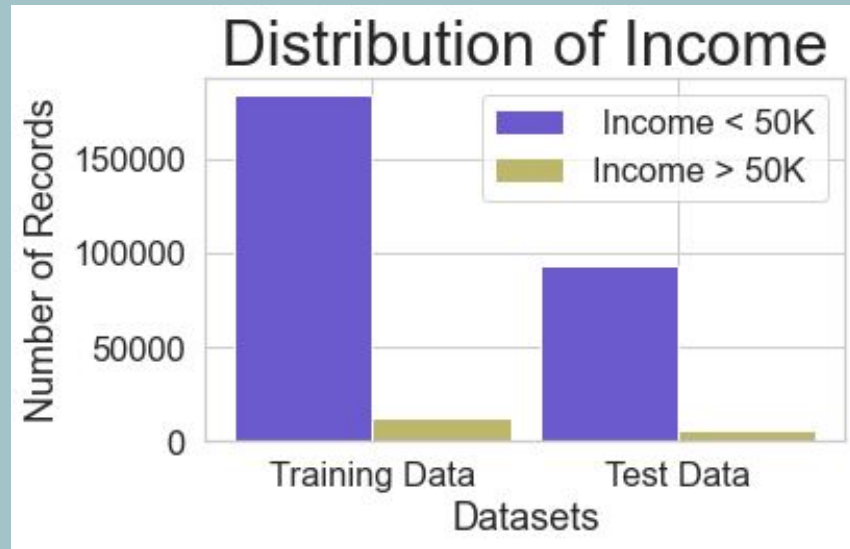
John Enright

Approach

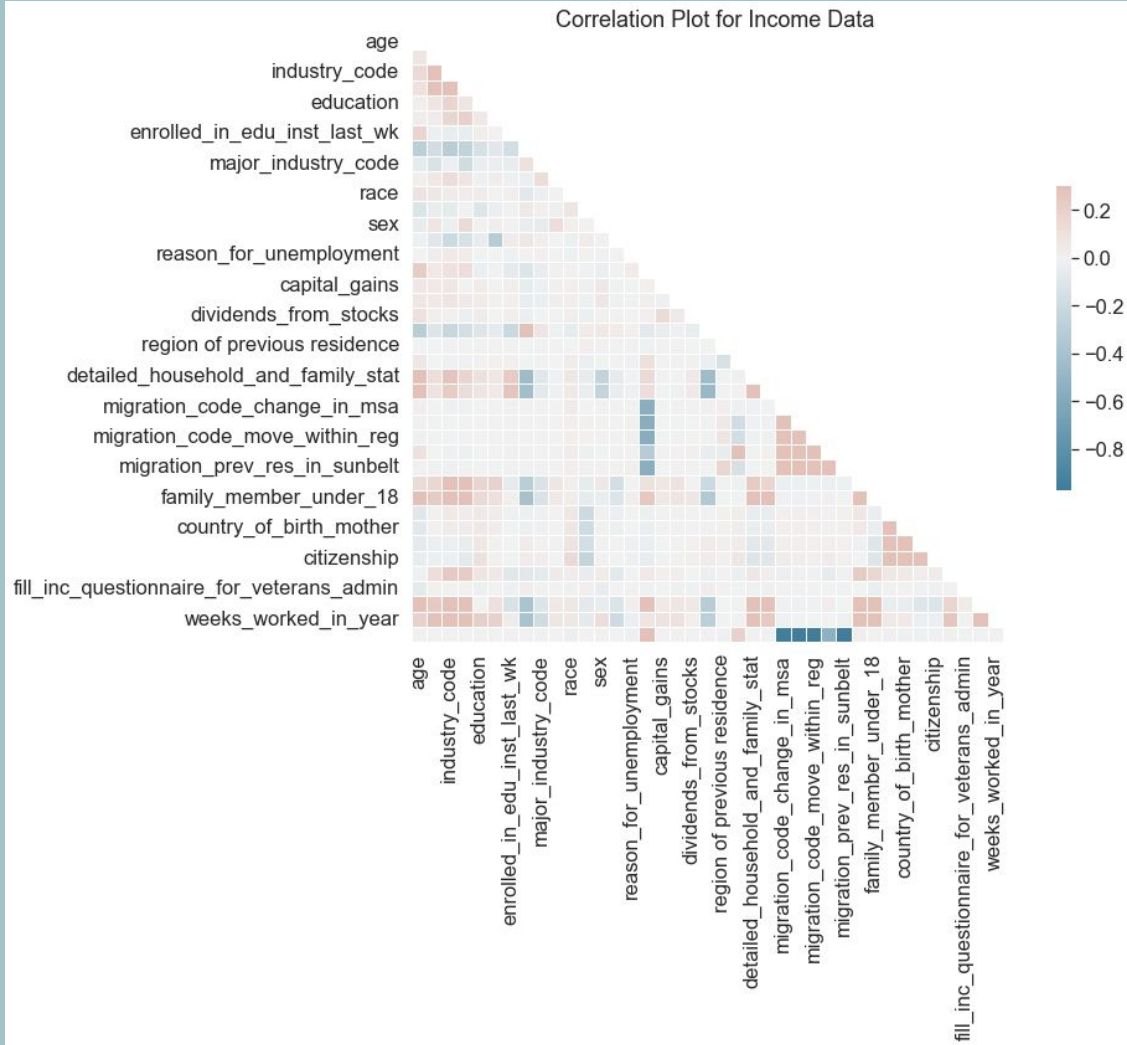


Exploratory Data Analysis

- Imbalanced data
 - 93% individuals make over \$50,000
 - Can introduce classification problems
- Duplicate or conflicting instances
 - Training → 46716
 - Test → 99762
- 40 characteristics (features) for an individual
 - 33 are continuous values
 - 7 are categorical
- Example features
 - Age
 - Sex
 - Race
 - Education level
 - Class of work

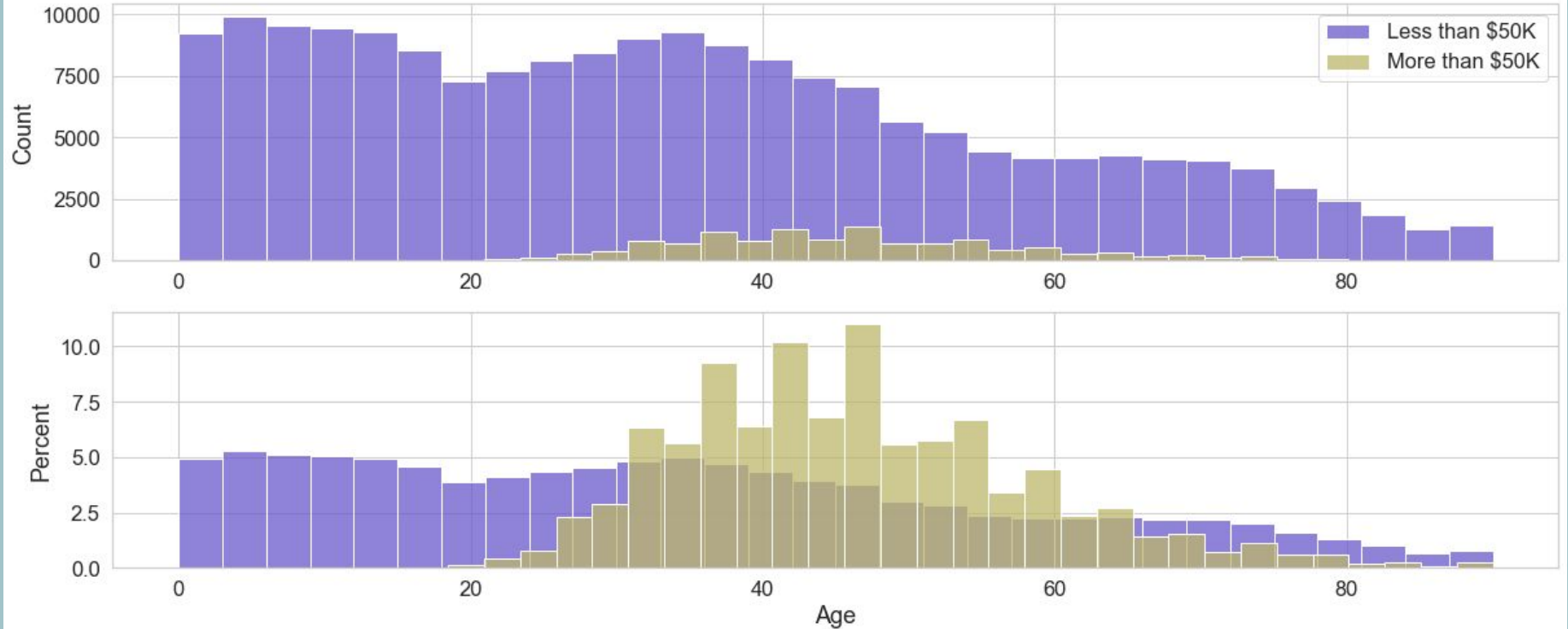


- Relationship

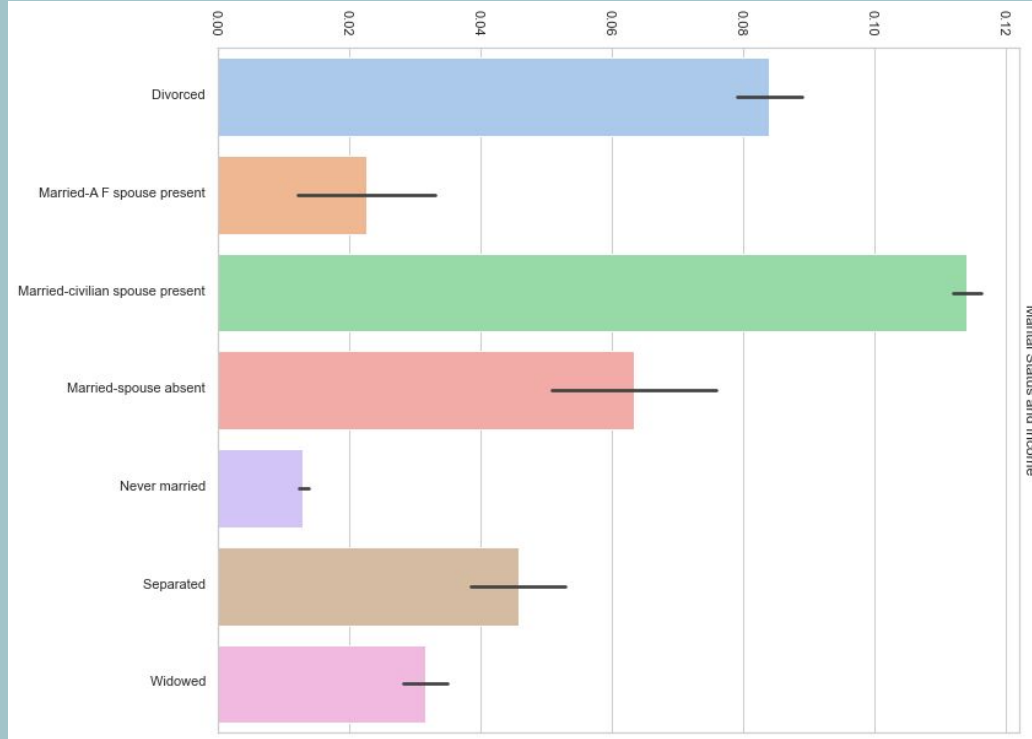


Exploratory Data Analysis

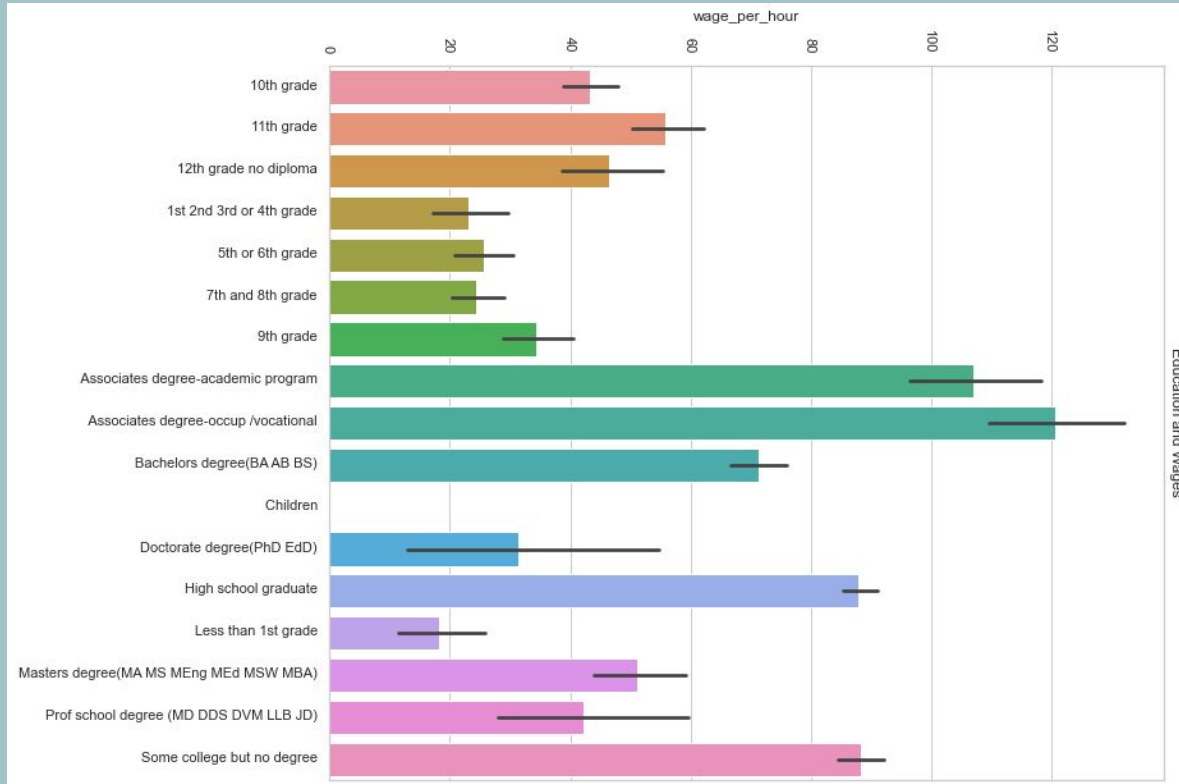
Distribution of Age by Income



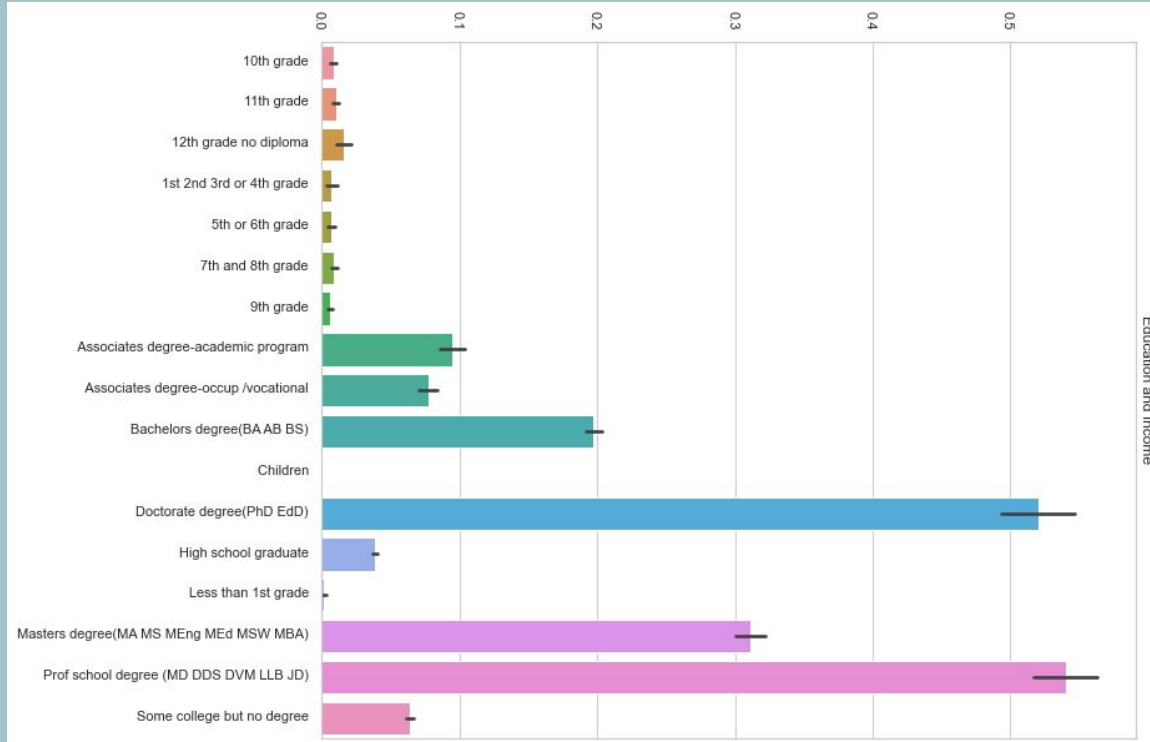
Exploratory Data Analysis - Marital Status



Exploratory Data Analysis - Education

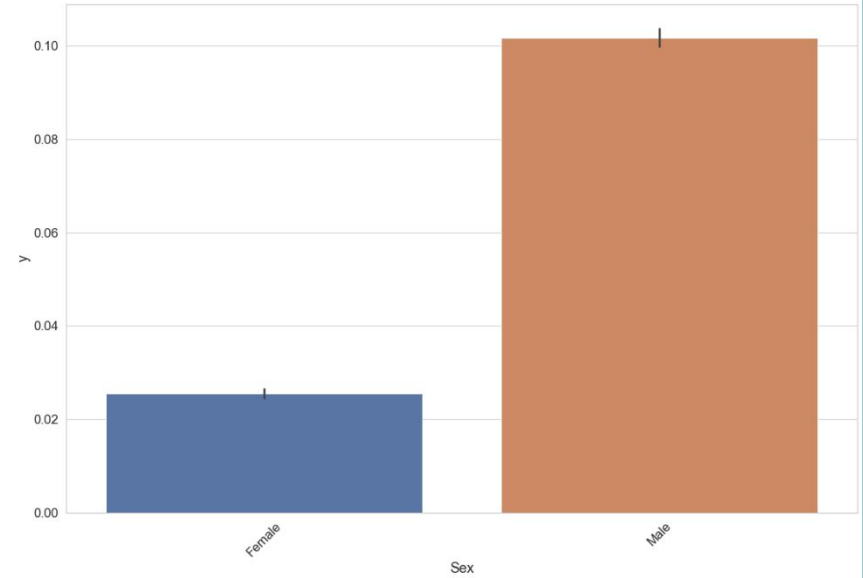
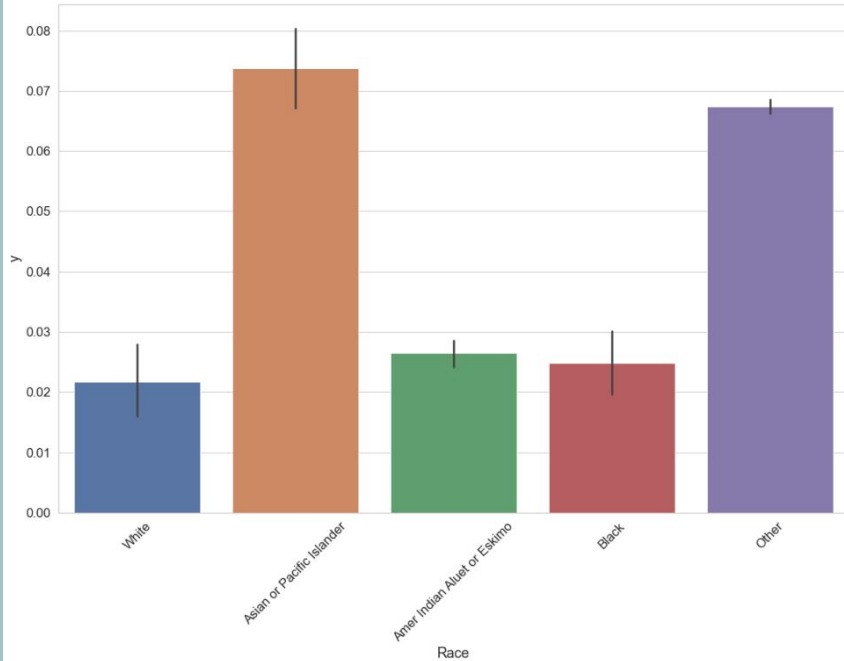


Exploratory Data Analysis - Education

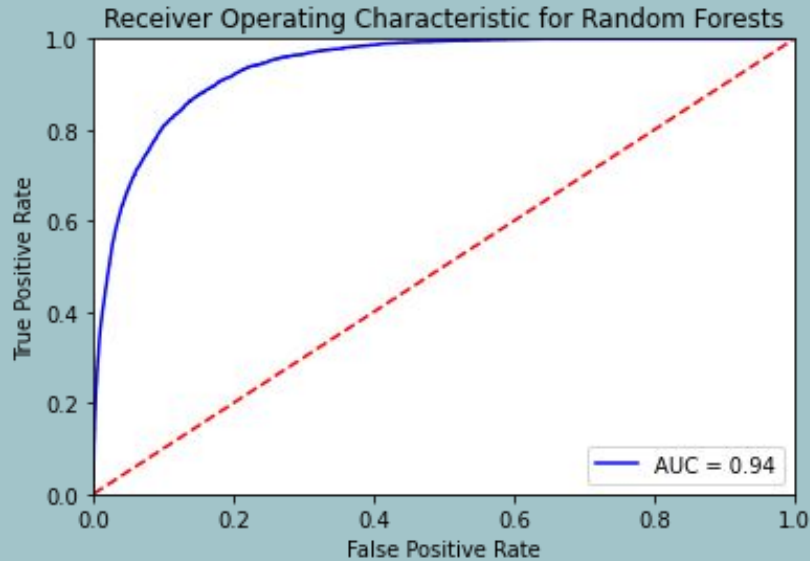


Exploratory Data Analysis - Race & Sex

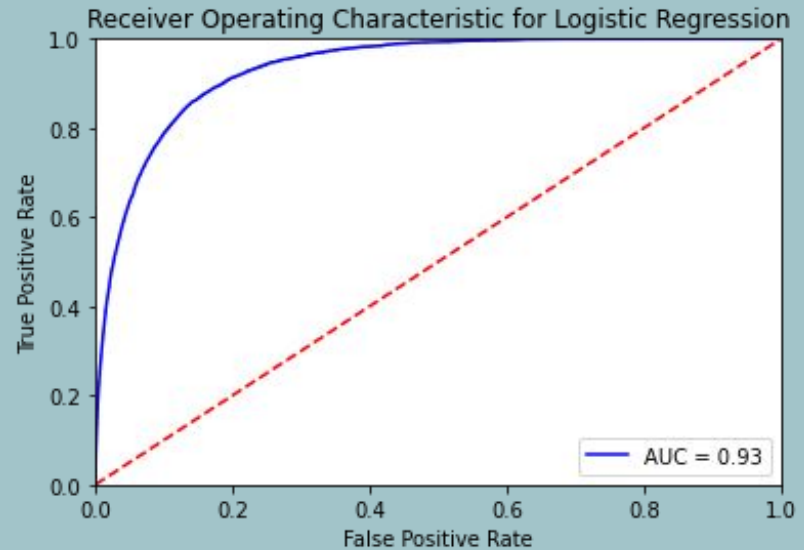
Individual Traits and Income



Model results - first pass

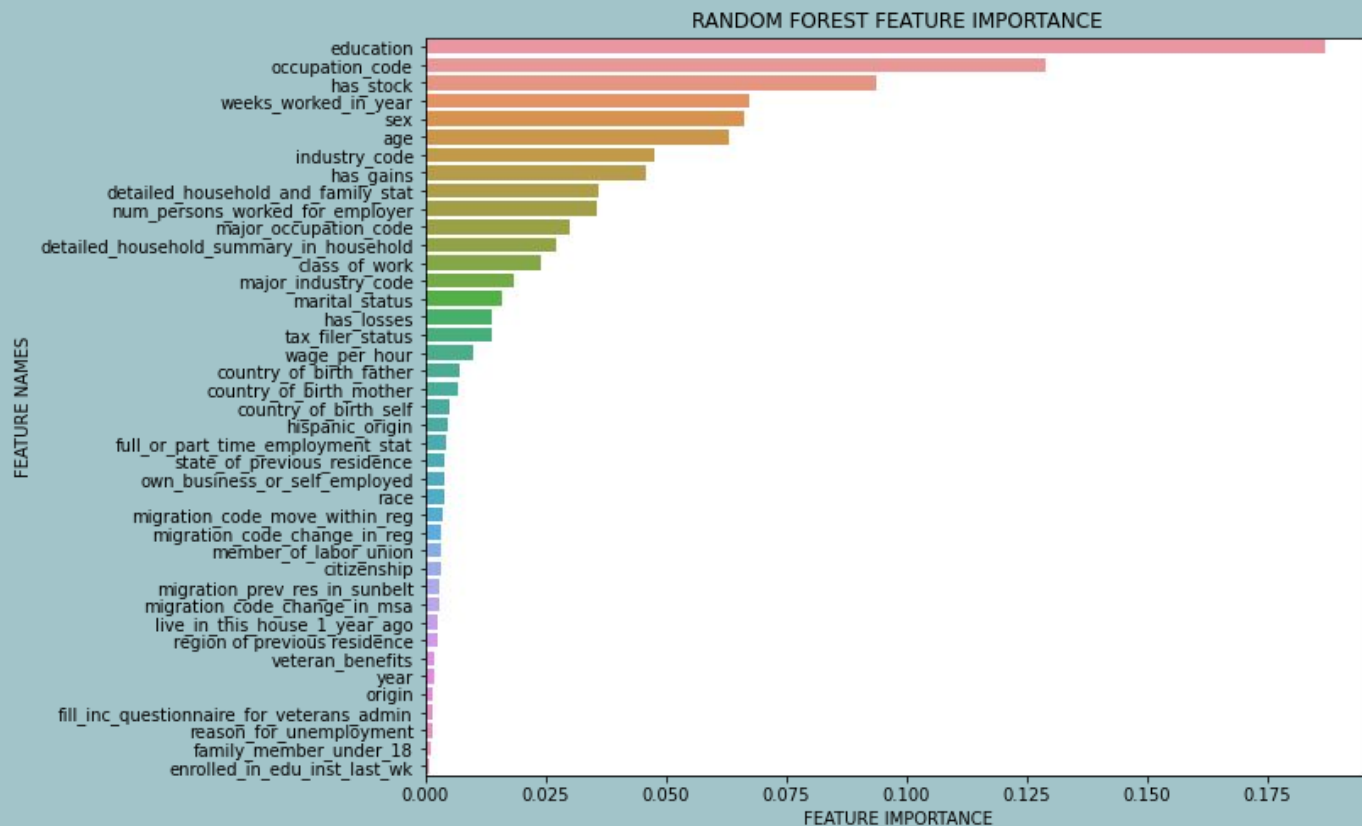


	F1 Score
Under \$ 50,000	0.97
Over \$ 50,000	0.41

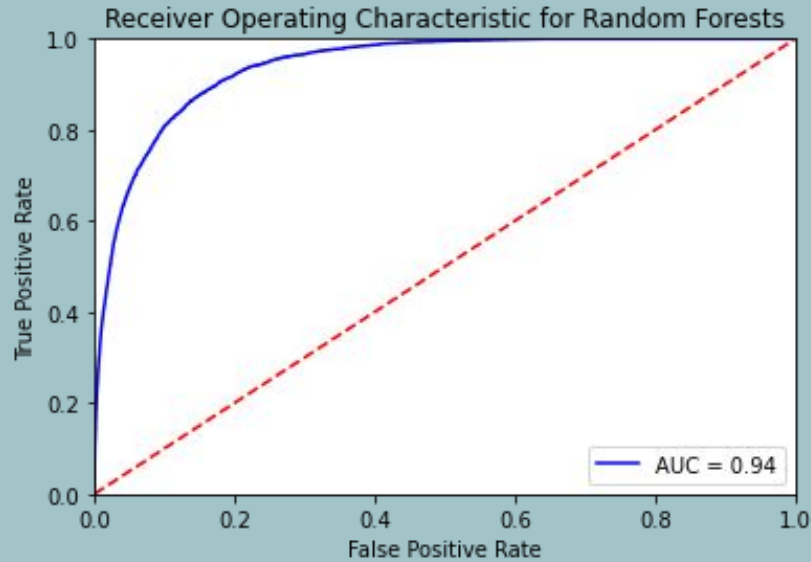


	F1 Score
Under \$ 50,000	0.97
Over \$ 50,000	0.44

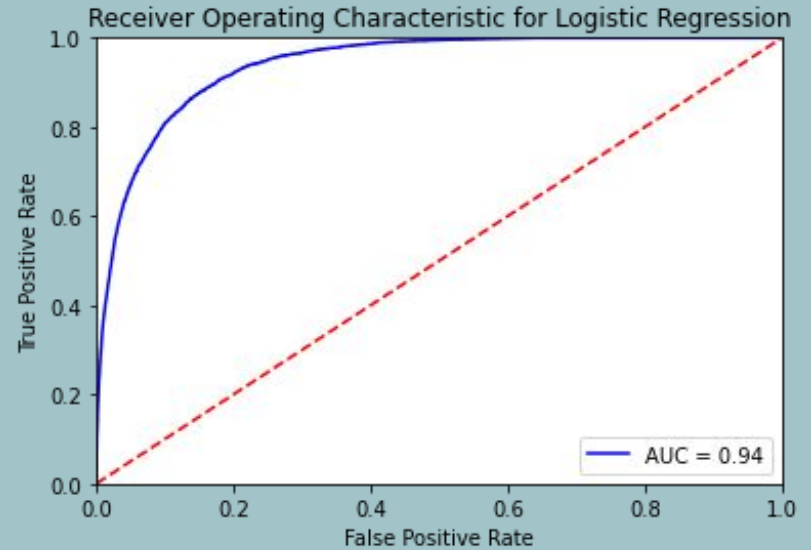
Model results - important features



Model results - reduced features



	F1 Score
Under \$ 50,000	0.97
Over \$ 50,000	0.46



	F1 Score
Under \$ 50,000	0.97
Over \$ 50,000	0.43

Thanks!

John Enright