

# PREDICTING INCOME LEVELS FROM U.S. CENSUS DATA

Michelle A. Caler

April 21, 2021

# OVERVIEW

- Introduction
- Associations with Income
- Models
- Results
- Conclusions

# INTRODUCTION

# MOTIVATION AND CENTRAL QUESTION

- Motivation: carry out first fully-independent data science project since earning Data Science Career Path certificate from Codecademy
- Central Question: Can I determine whether a person makes over \$50K a year?

## THE DATA: WHERE AND WHY

- Income and demographics data compiled from 1994 Census Bureau database
- Downloaded from [this Kaggle repository](#)
  - Repository owner: UCI Machine Learning group
- Public Domain data; compiled by UCI Machine Learning group so a clear answer to the question likely to exist
- Modestly-sized data set with a mix of quantitative and qualitative features

## THE DATA: WHAT AND HOW

- Data extracted from 1994 Census Bureau database
  - Paper cited by UCI Machine Learning Group: Ron Kohavi, "[Scaling Up the Accuracy of Naive-Bayes Classifiers: a Decision-Tree Hybrid](#)", Proceedings of the Second International Conference on Knowledge Discovery and Data Mining, 1996
- 32,561 records, 15 features
  - No null values!

# THE DATA: FEATURES

Quantitative	Qualitative
age	workclass
fnlwgt	education
education.num	marital.status
capital.gain	occupation
capital.loss	relationship
hours.per.week	race
	sex*
	native.country
	income*

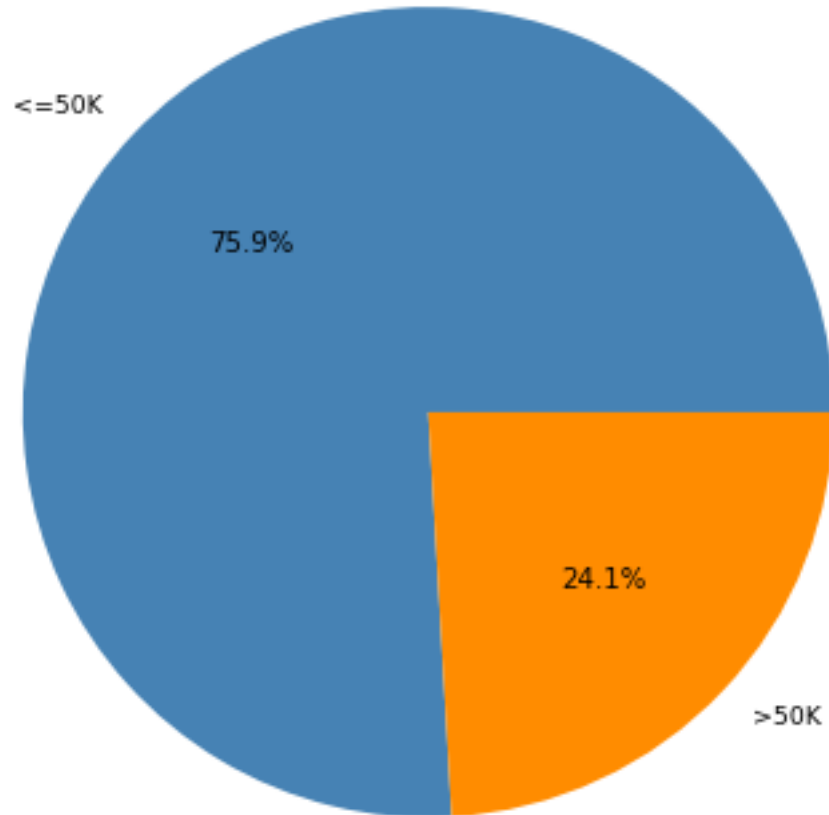
\*: binary categorical variable

# ASSOCIATIONS WITH INCOME



# INCOME CATEGORIES

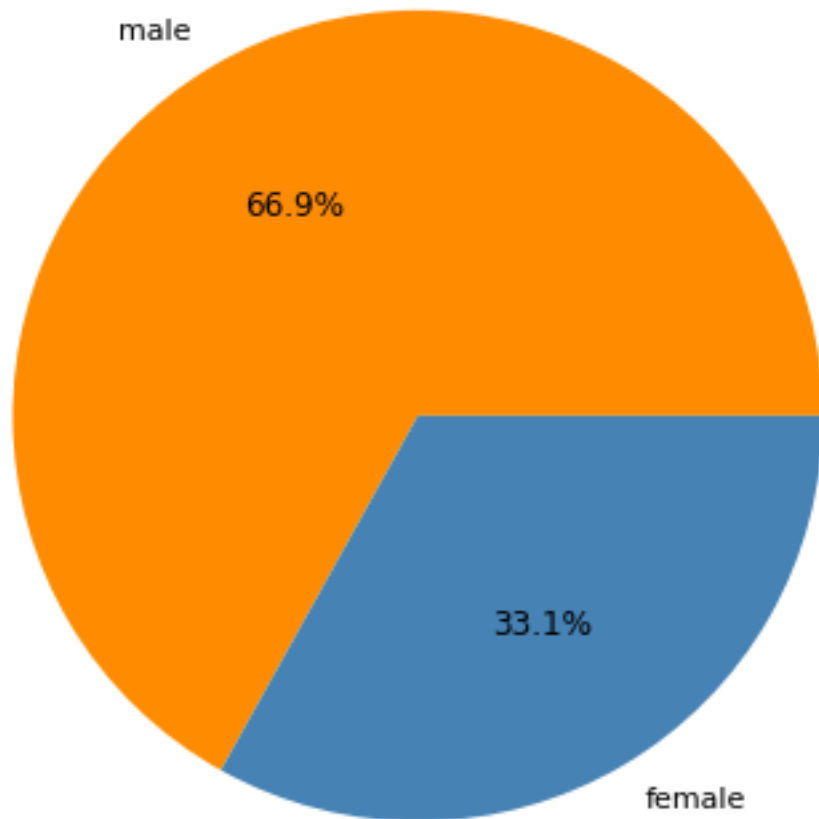
Percentage Breakdown of Income Categories



- ◇ Low-income: <= 50K
  - ◇ Roughly  $\frac{3}{4}$  of records
- ◇ High-income: > 50K
  - ◇ Roughly  $\frac{1}{4}$  of records

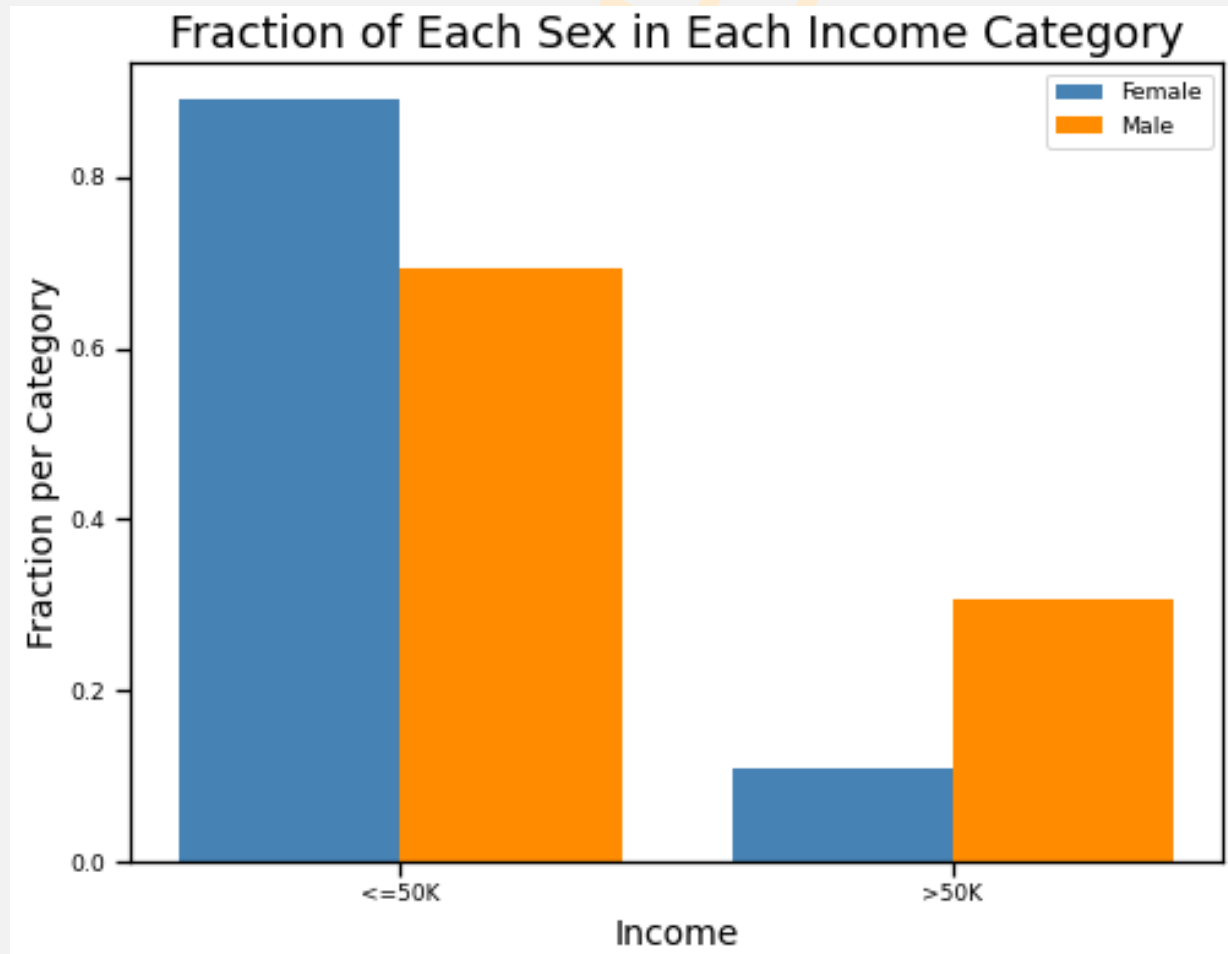
# SEX AND INCOME

Percentage Breakdown of Biological Sex



- ◇ Sex = Biological Sex
- ◇ Roughly 2/3 of records are males
- ◇ Roughly 1/3 of records are females

# SEX AND INCOME



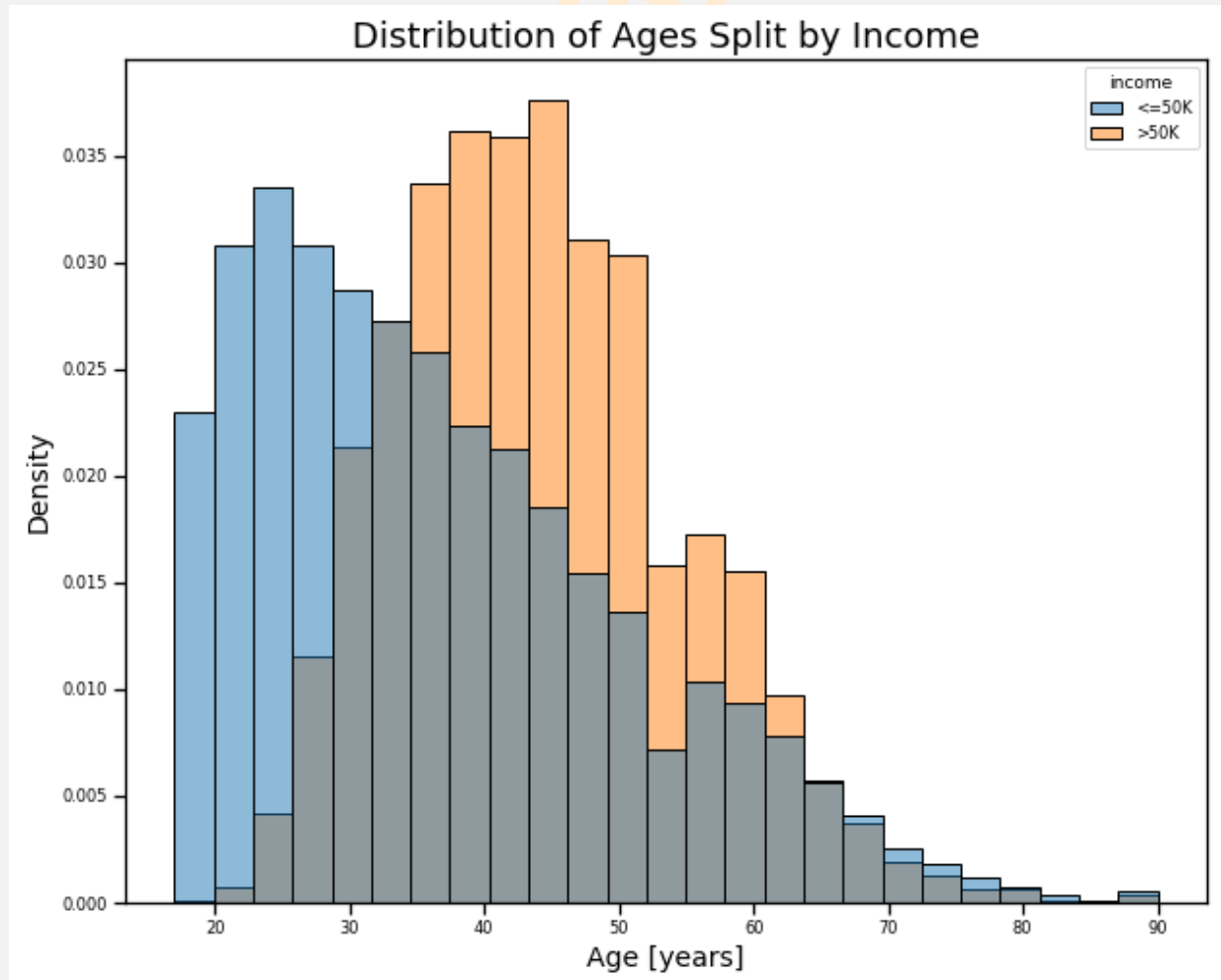
◇ Females:

- ◇ Roughly 85% in low-income category
- ◇ Roughly 15% in high-income category

◇ Males:

- ◇ Roughly 65% in low-income category
- ◇ Roughly 35% in high-income category

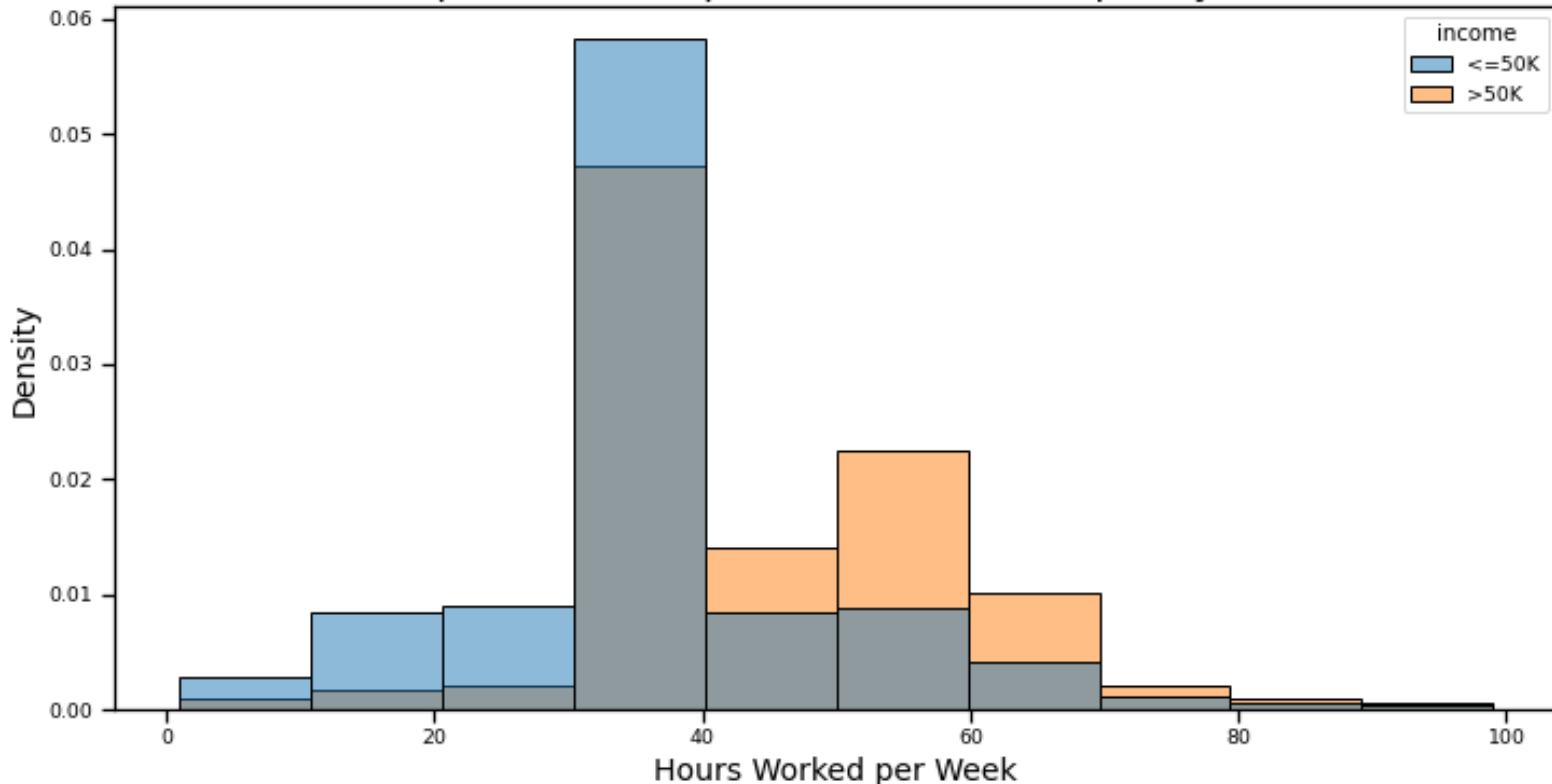
# AGE AND INCOME



- ◇ Low-income and high-income distributions have different shapes
  - ◇ Low-income: gradual decay from peak
  - ◇ High-income: skewed bell curve
- ◇ Low-income and high-income distributions peak at different ages
  - ◇ Low-income: peak around 25 yrs. old
  - ◇ High-income: peak around 45 yrs. old

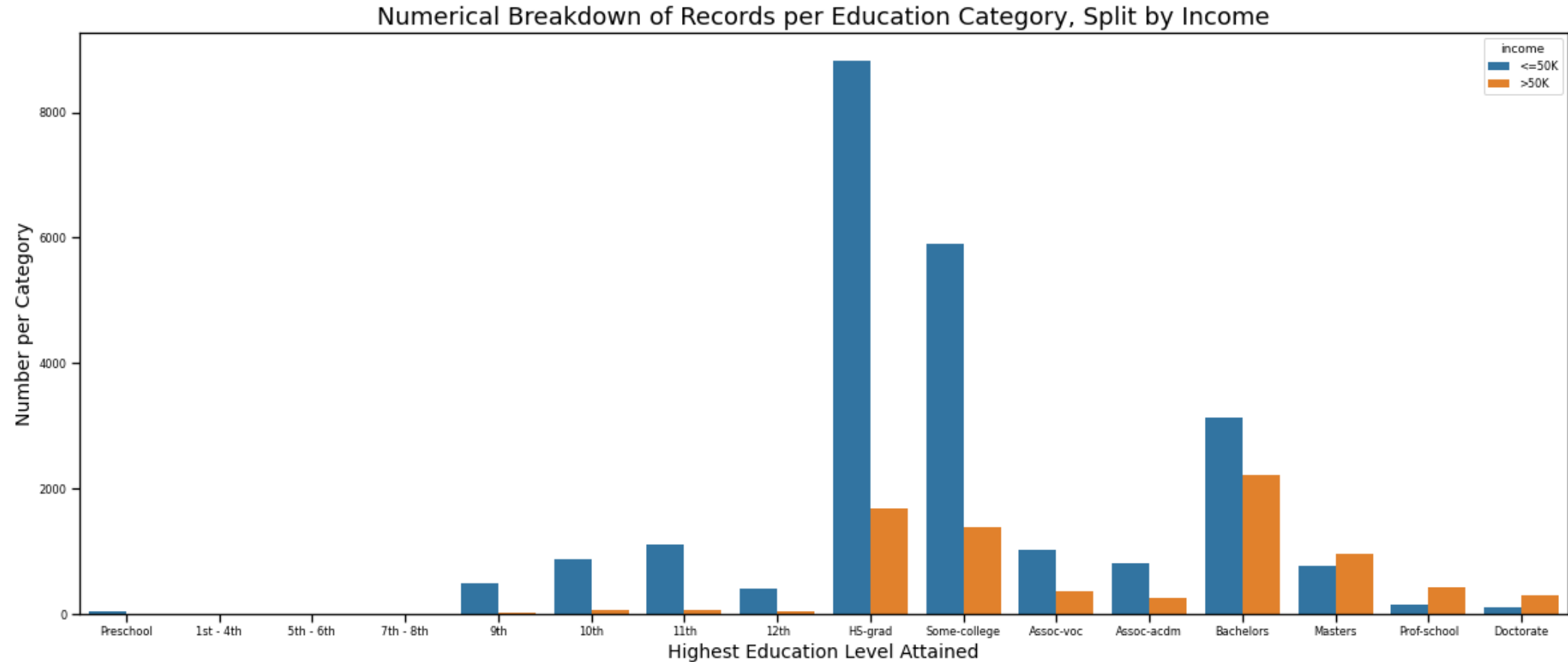
# HOURS WORKED PER WEEK AND INCOME

Self-Reported Hours per Week Worked, Split by Income



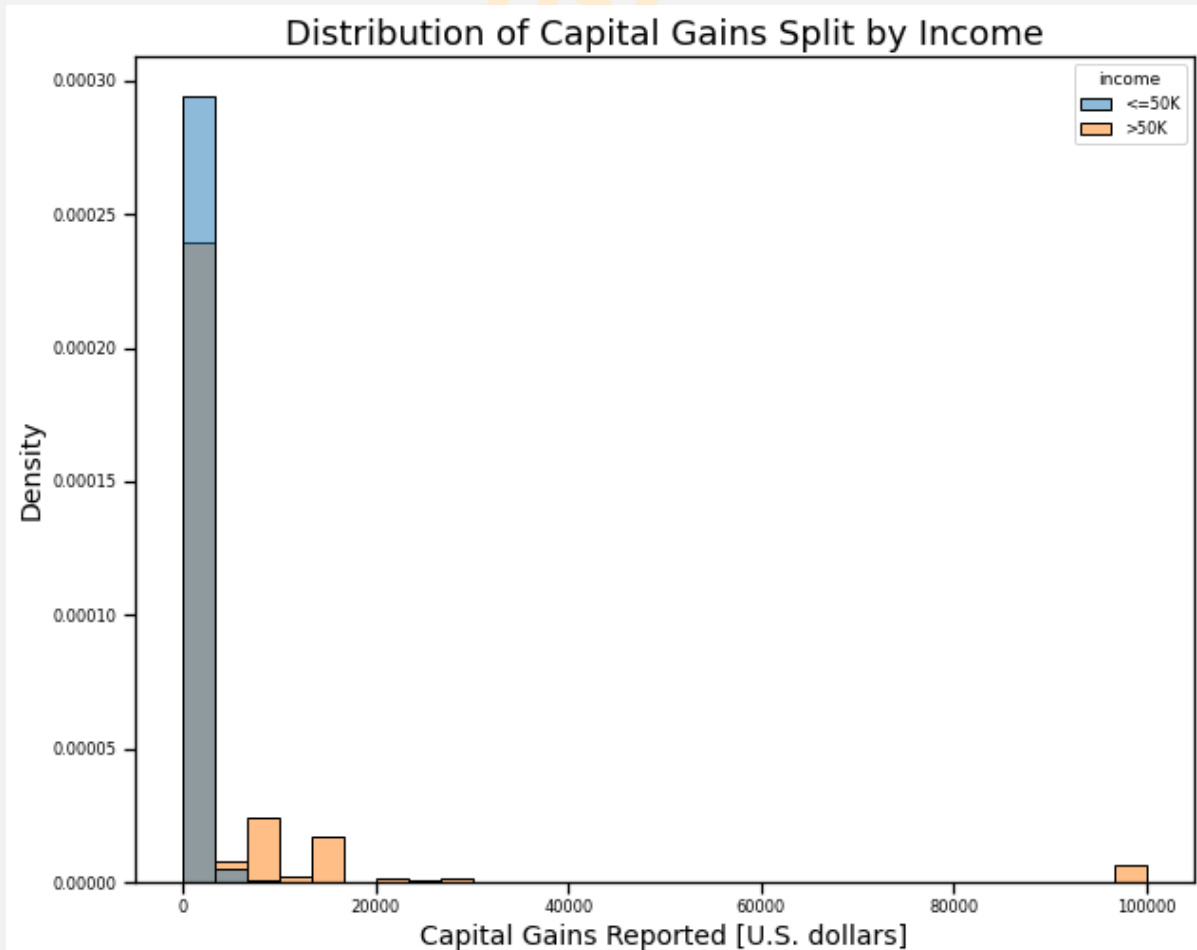
- ◇ Low-income: peaks at 35-40 hours per week
  - ◇ Plateau from 15-35 hours per week
  - ◇ Plateau from 40-60 hours pre week
- ◇ High-income: peaks at 35-40 hours per week and 55-60 hours per week
  - ◇ Peak at 35-40 hours per week much sharper
  - ◇ Skew left

# EDUCATION AND INCOME



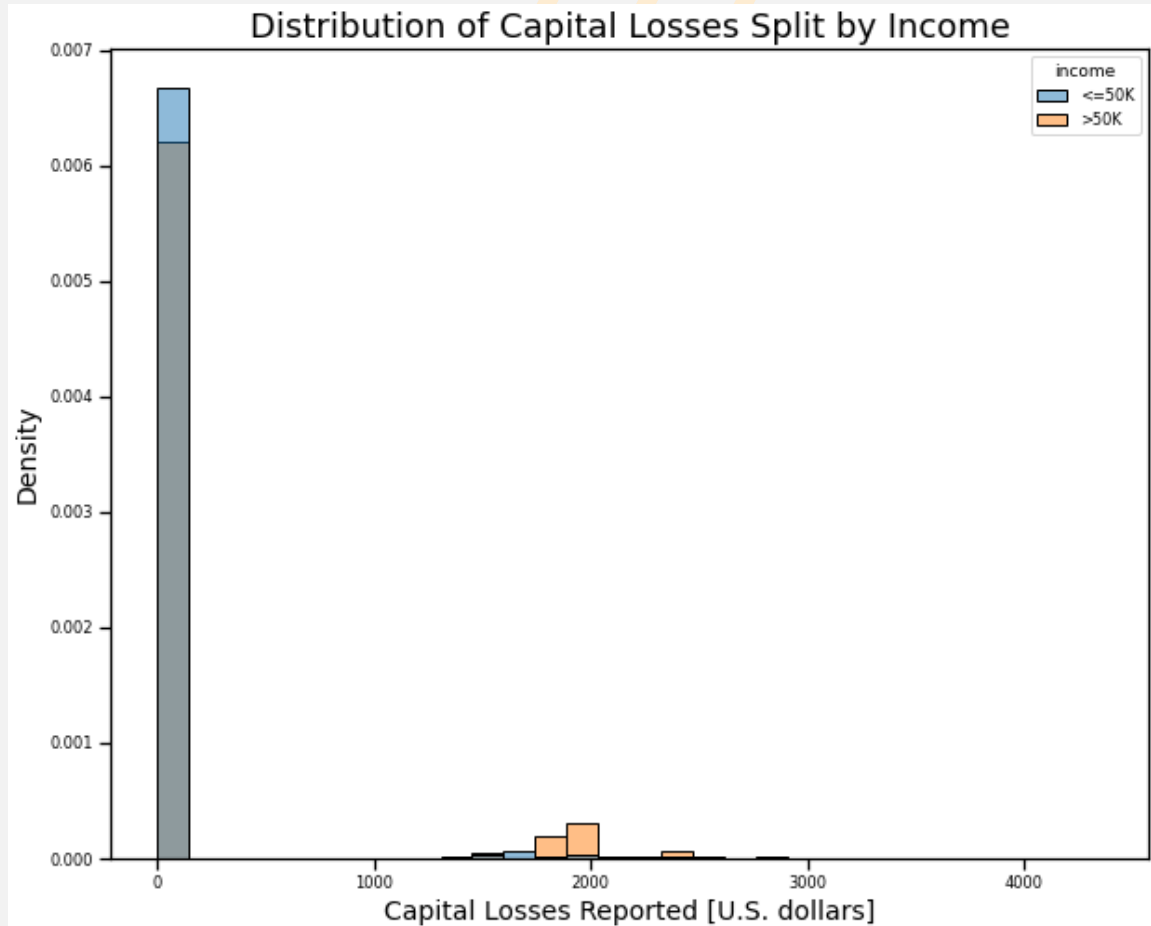
◇ General trend: the higher the education level attained, the higher the fraction of high-income records

# CAPITAL GAINS AND INCOME



- ◇ Low-income: Little to no capital gains
- ◇ High-income: Mostly little to no capital gains, but significant minority with gains > 10K

# CAPITAL LOSSES AND INCOME



- ◇ Low-income: Little to no capital losses
- ◇ High-income: Mostly little to no capital gains, but significant minority with losses > 15K



# SUPERVISED MACHINE LEARNING MODELS

## MODEL INPUTS

- Features used to train models:
  - Sex, converted to integers (1: Female, 0: Male)
  - Education, converted to integers (0: Preschool through 16: Doctorate)
  - Age
  - Hours Worked per Week
  - Capital gains
  - Capital losses
- Prediction labels:
  - 0: low-income record
  - 1: high-income record

## MODEL INPUTS

- 75% of data used to train models  
25% of data used to test/validate models
- Training and test data scaled such that all features lay on same scale
- Exact same training data and test data used for all models

# MODELS TRAINED TO FEATURES

(AND USED TO PREDICT INCOME LEVEL)

- Logistic Regression Classifier
- Support Vector Machine Classifier
- K-Nearest Neighbours Classifier
- Single Decision Tree Classifier
- Random Forest Classifier

# RESULTS

## MODEL PRECISION, ACCURACY, RECALL

**Highest  
precision,  
highest  
accuracy**



Model	Precision	Accuracy	Recall
Logistic Regression	71.7%	81.8%	41.0%
Support Vector Machine	76.7%	83.0%	42.9%
K-Nearest Neighbours	72.9%	82.2%	41.9%
Single Decision Tree	66.9%	81.9%	49.9%
Random Forest	69.3%	82.8%	51.7%

**Best  
balances  
precision  
and recall**



# CONCLUSIONS

## CONCLUSIONS

- I was able to build a supervised machine learning model to predict income level with good precision and accuracy
- Best model: Support Vector Machine classifier
  - 76.7% precision
  - 83.0% accuracy
  - 42.9% recall
- Best Balanced model: Random Forest classifier
  - 69.3% precision
  - 82.8% accuracy
  - 51.7% recall



## IMPROVEMENTS

- Use a machine learning classification model which accepts qualitative as well as quantitative data
  - Nominal categorical variables with clear associations with income unused
- Explore  $\gamma/C$  parameter space of Support Vector Machine classifier
  - Is current model a local max or a global max?

## CITATIONS

- Ron Kohavi, "[Scaling Up the Accuracy of Naive-Bayes Classifiers: a Decision-Tree Hybrid](#)", Proceedings of the Second International Conference on Knowledge Discovery and Data Mining, 1996.
- The [scikit-learn Python machine learning library](#): Machine Learning in Python, Pedregosa et al., JMLR 12, pp. 2825-2830, 2011
- After completing my project, I came across [this](#) stellar write-up from (then) UCSD students Chet Lemon, Chris Zelazo, and Kesav Mulakaluri. Their treatment of the problem no doubt influenced my own write-up of my work.