

Data Scientist Technical Assessment and Presentation

How to identify a \$50k+ a year income earner?

David Rickheim | 24JUL2024

About me - David Rickheim



Agenda

1. Problem Statement
2. Exploratory Data Analysis
3. Data Preparation
4. Data Modeling
5. Model Assessment
6. Results

1: The problem

Problem statement

Can we predict who makes over \$50k annually?

How may we differentiate this group (what predicts a >\$50k earner)?

Importance

Gov't and company policy decision making

Determines qualifications for many things, such as:

- SNAP
- School lunches
- Medicaid

Census Data

Annual Social and Economic (ASEC) Supplement

- 1993 / 1994

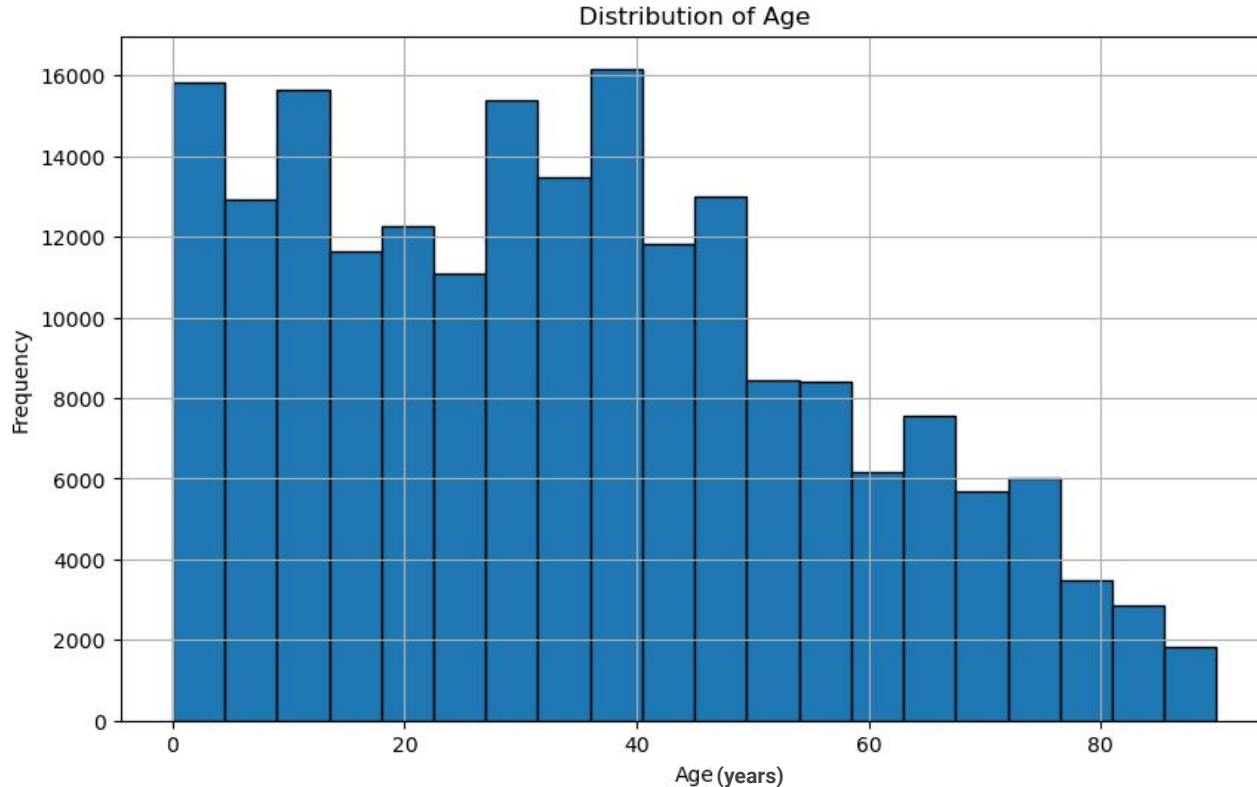
Hourly wage, hours worked, demographic info, occupation, education

2. Exploratory Data Analysis

The data - key statistics

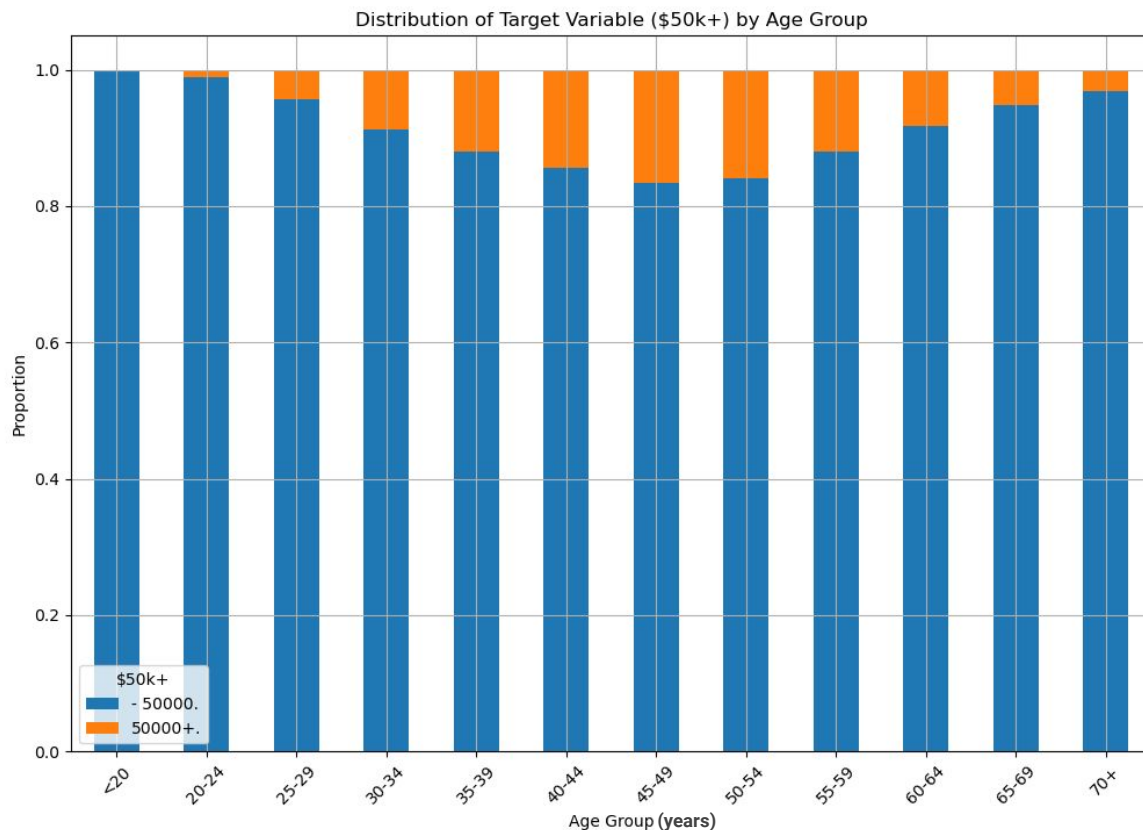
	Training Data	Testing Data
Inputs / features / columns	41	41
Target variable	Income bracket (\$50k+)	Income bracket (\$50k+)
Rows / Records	199,522	99,761
# of \$50k+	12,382	6,186
% of \$50k+	6.2%	6.2%

Ages of records



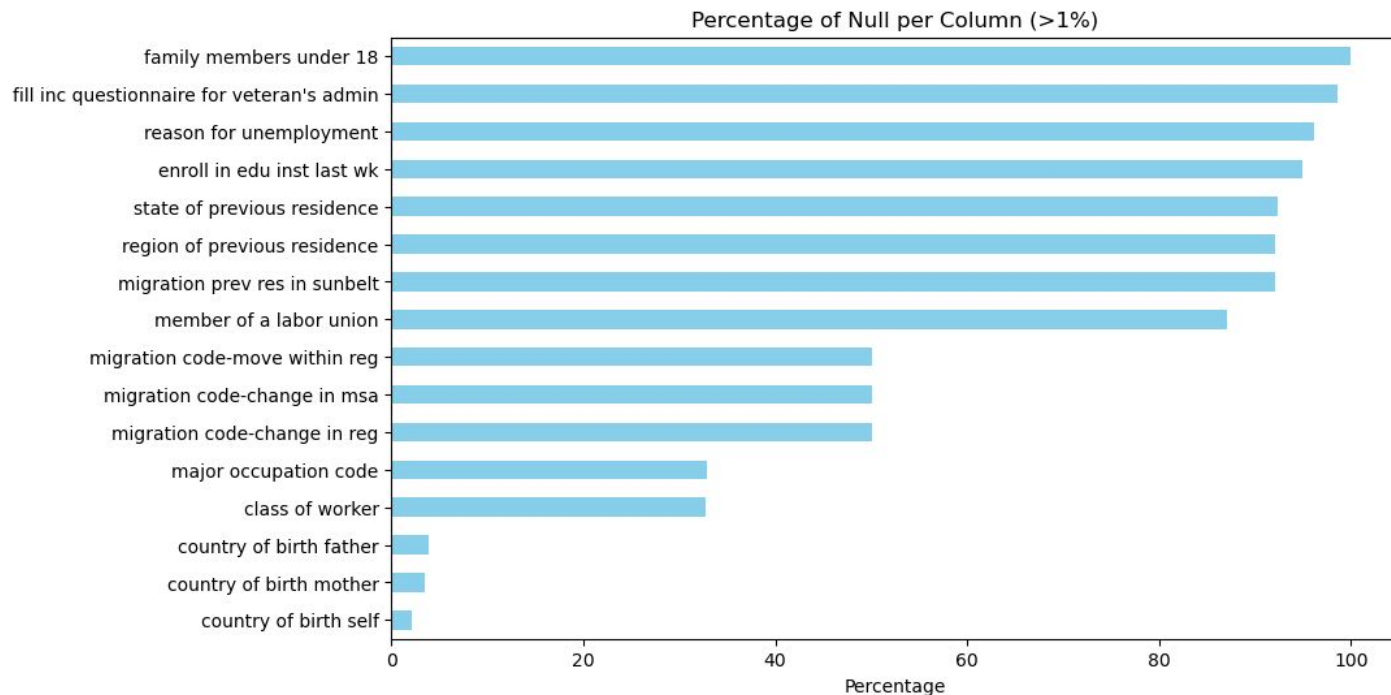
- Mean: 34.5 years
- Median: 33 years

Income distribution by age group



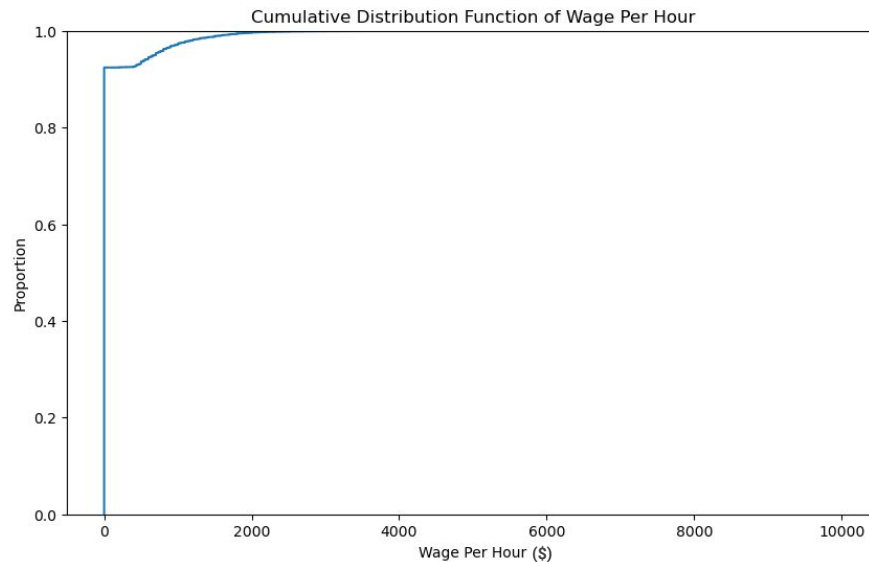
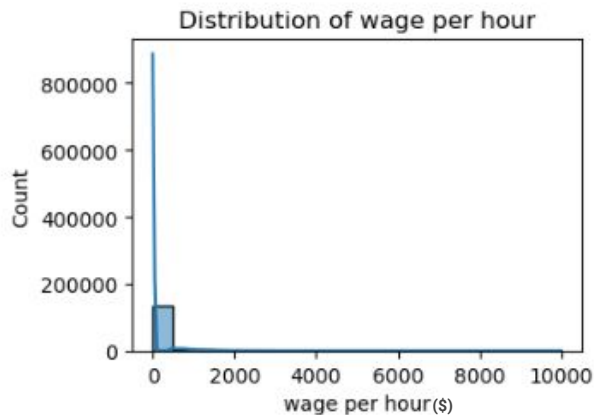
- Children very rarely make \$50k
- Vast majority of juveniles don't have a job code or industry code
- Not to be included in analysis

Not in universe - a common code for “nulls”



Question mark (“?”) a common null value for country fields

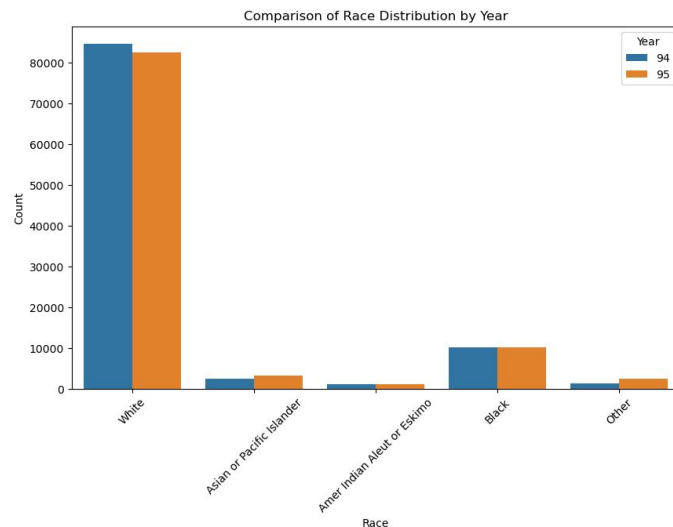
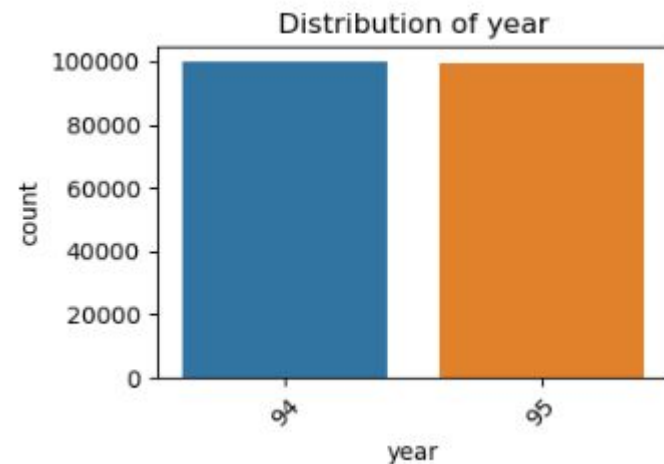
Skewness of wages



- Wage alone **cannot** identify all high income earners since >90% are empty
- Other income sources (capital gains, dividends) also have a high number of "\$0" entries

We have two years of data

- Possibly describes the same populations
- Similarity between the data represented by years '94' and '95' are analyzed using statistical tests
 - Chi-squared & t-tests
- No significant differences
- Extra attention given to protected classes
 - Race, sex, hispanic origin, citizenship, veterans benefits



3. Data Preparation

Data Preparation Challenges

Nulls & Skewness

Drop eleven columns that are >50% null

Impute nulls in other columns

Scale / transform skewed variables

Outlier check

Year column

Difficult to interpret as data could be longitudinal

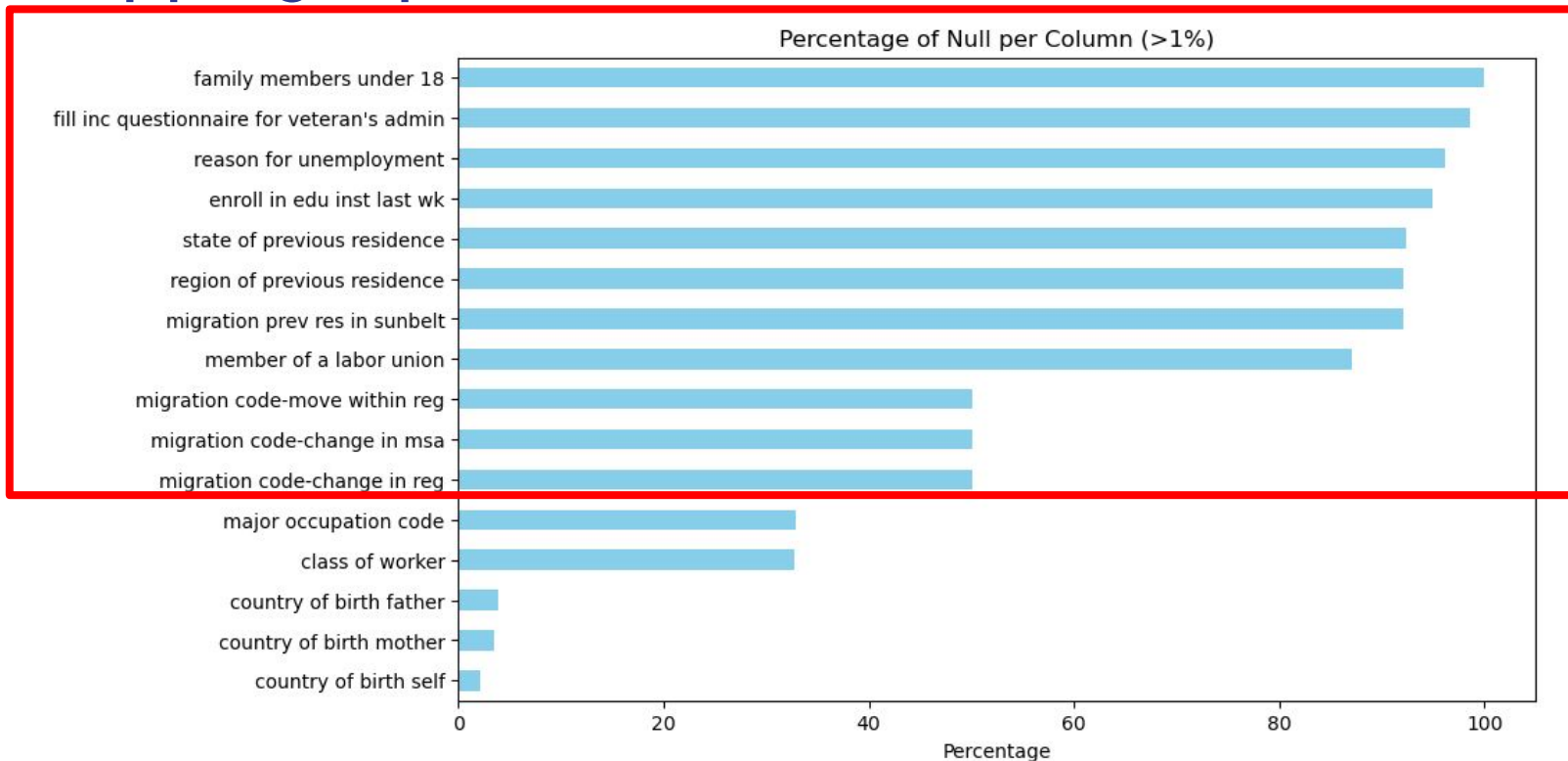
We elect to choose one year ('94') as the data seems to be very similar

Sample selection

Who do we want to predict?

Children <18 removed

Dropping inputs with >50% null



Feature engineering: annual salary

- Noticeably missing from the dataset is annual salary.
- Assuming 40 hour work weeks, we will use "wage per hour" and "number of weeks worked in year" to estimate annual salary.



4. Data Modeling

Model Selection

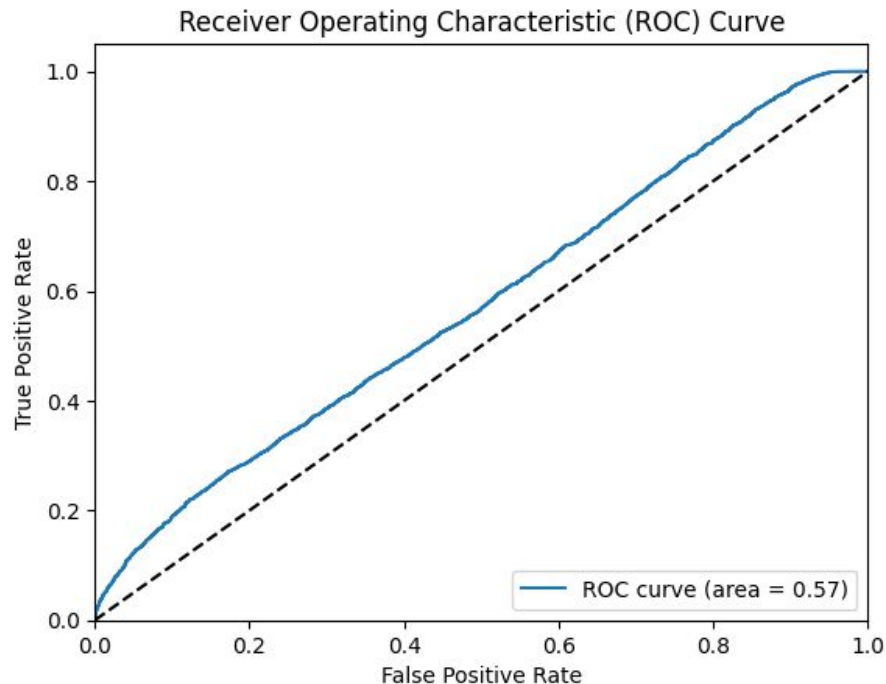
- Three models:
 - logistic regression - interpretable, relatively simple
 - random forest - interpretable, robust against outliers
 - XGBoost - superior performance, handles missing data better
- Main scoring metric will be AUC (area under curve)
 - Robust when dealing with an imbalanced dataset
 - Makes for easier comparisons between models
- Recall of >\$50k income earners also considered
 - We would like to avoid missing too many high earners

5. Model Assessment

Logistic Regression

	Precision	Recall
Over 50k	0.19	0.11
Under 50k	0.93	0.96

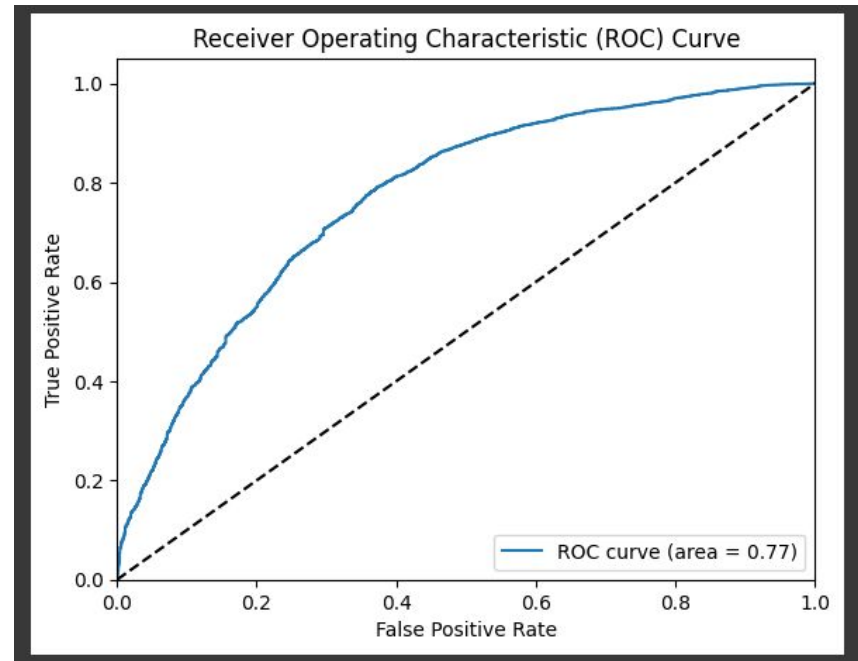
- AUC ROC: 0.57
- Model performs poorly
- 11% of >\$50k captured



Random Forest

	Precision	Recall
Over 50k	0.40	0.15
Under 50k	0.93	0.98

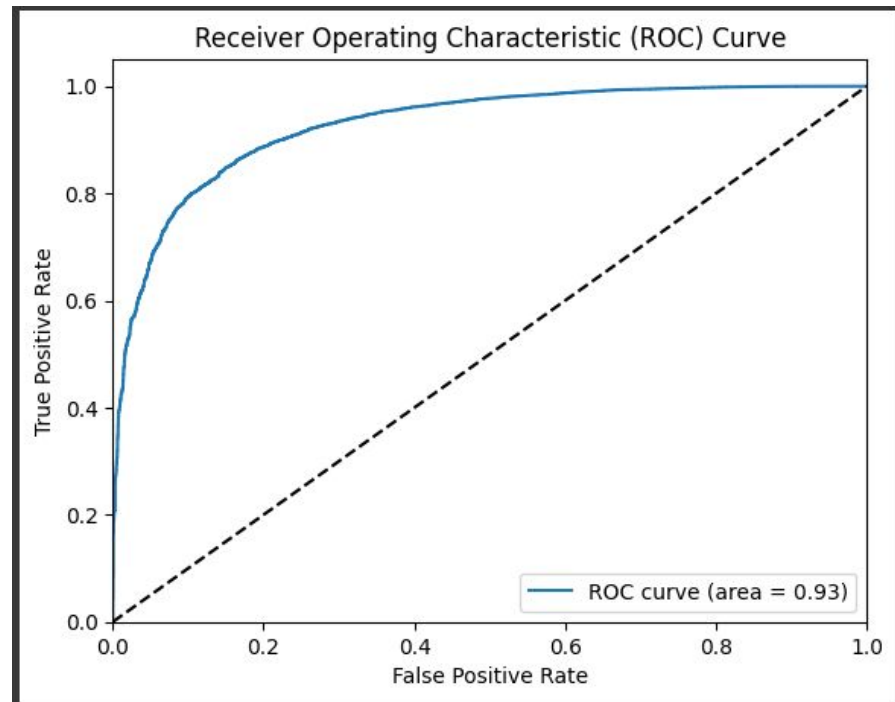
- AUC ROC: 0.77
- Model performance is fair
- 15% of >\$50k captured



XGBoost

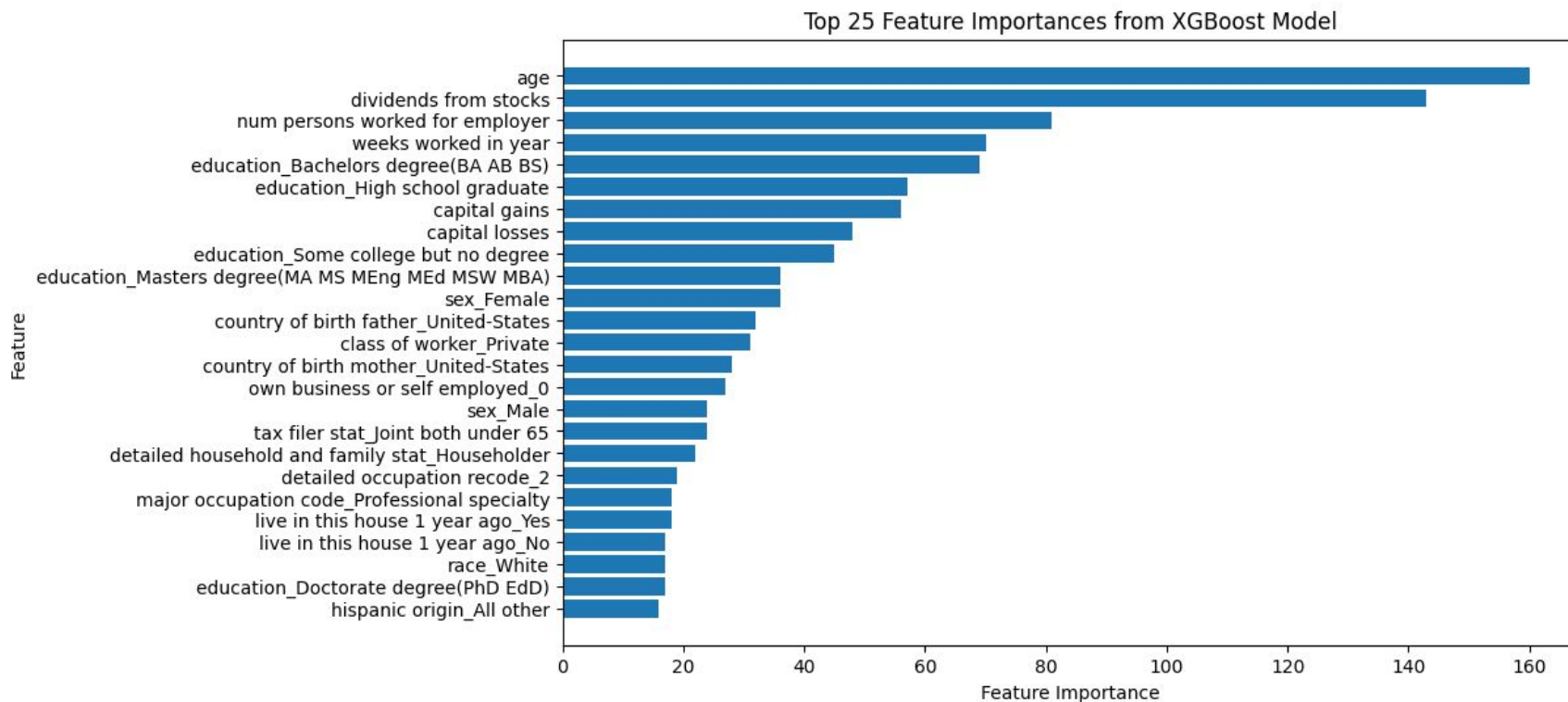
	Precision	Recall
Over 50k	0.64	0.52
Under 50k	0.96	0.97

- AUC ROC: 0.93
- Model performance fairly well
- 52% of >\$50k captured
- 64% of selected >\$50k are true



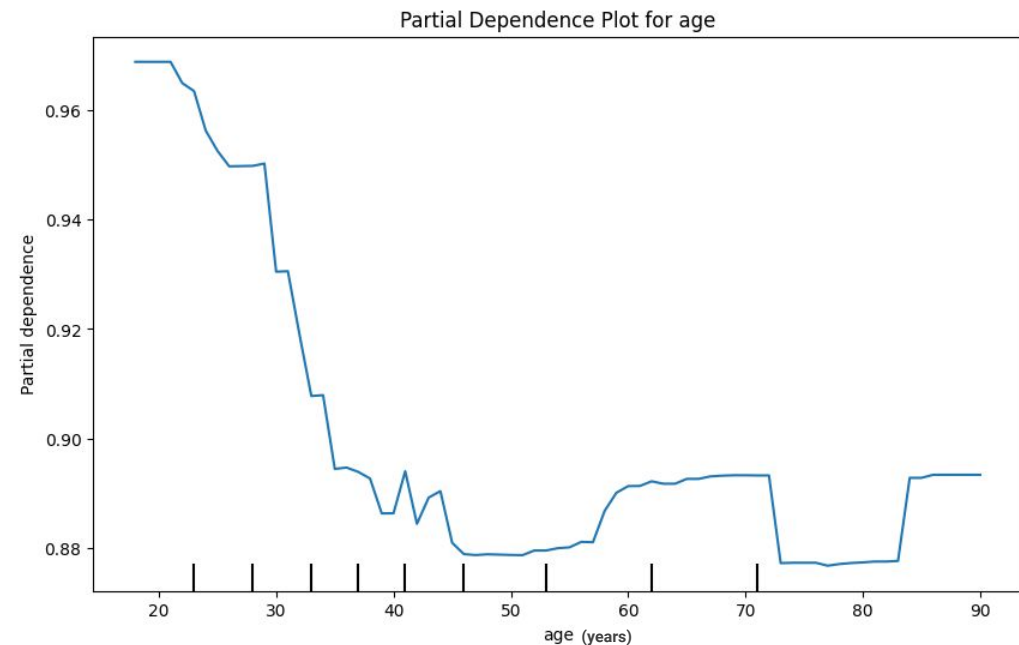
6. Results

Most important variables



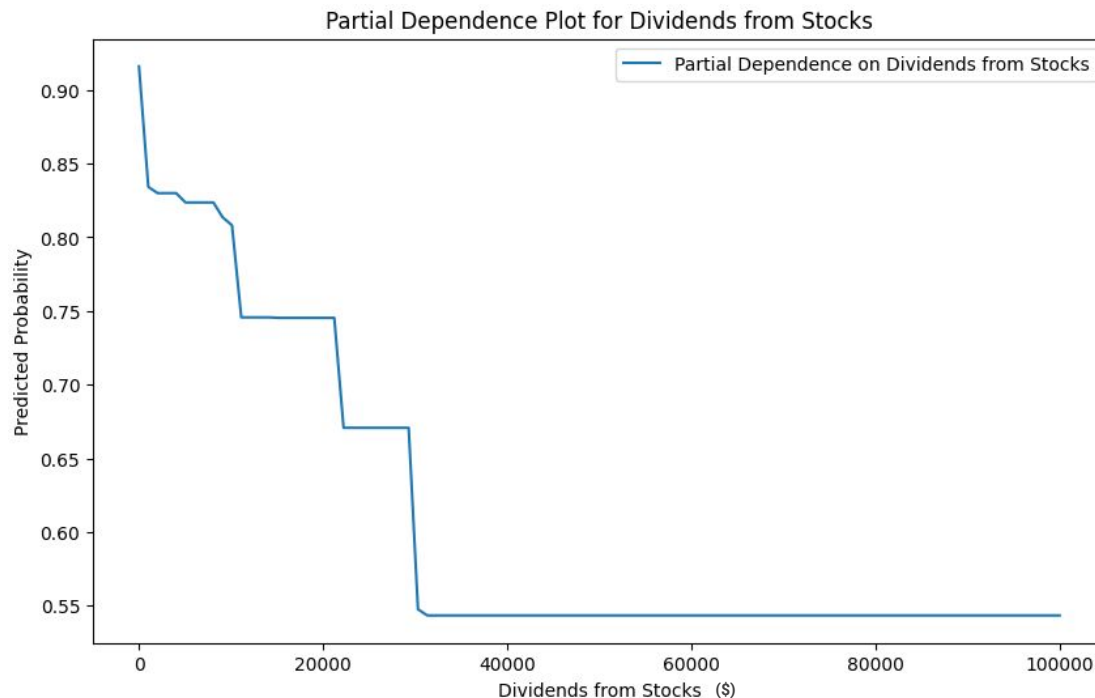
How age impacts income

- Higher plot line means more likely to make less than \$50k
- Likelihood of >\$50k **increases** as one approaches ~45 years old, then stagnates



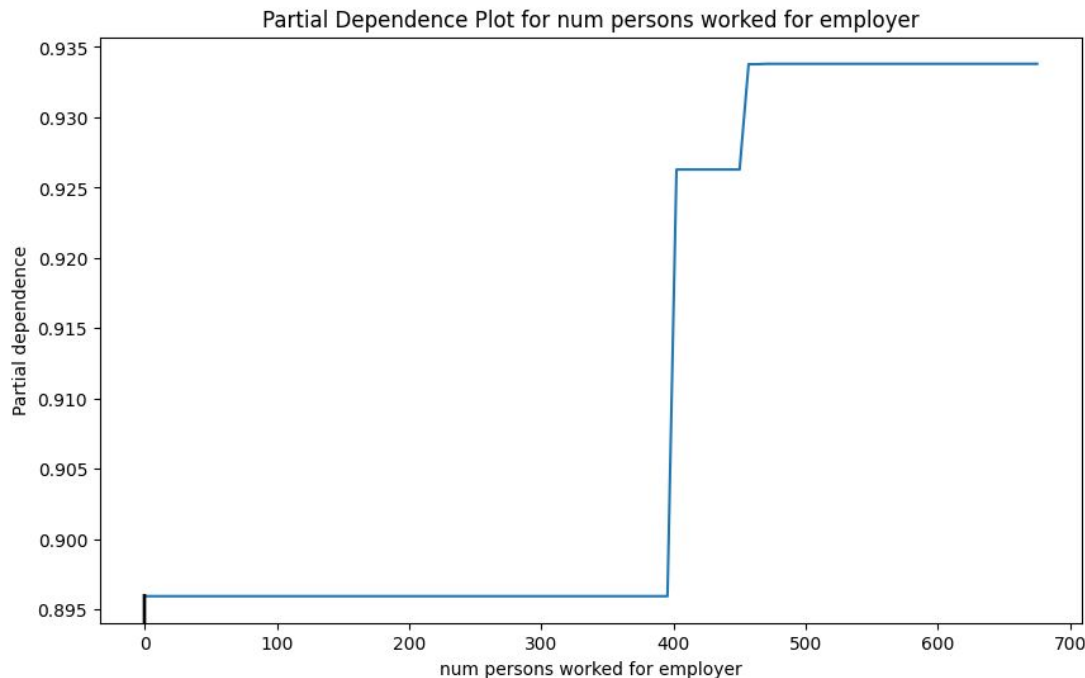
Dividend earners are likely to make >\$50k

- Likelihood of >\$50k **increases** as one approaches \$30k worth of stock dividends annually

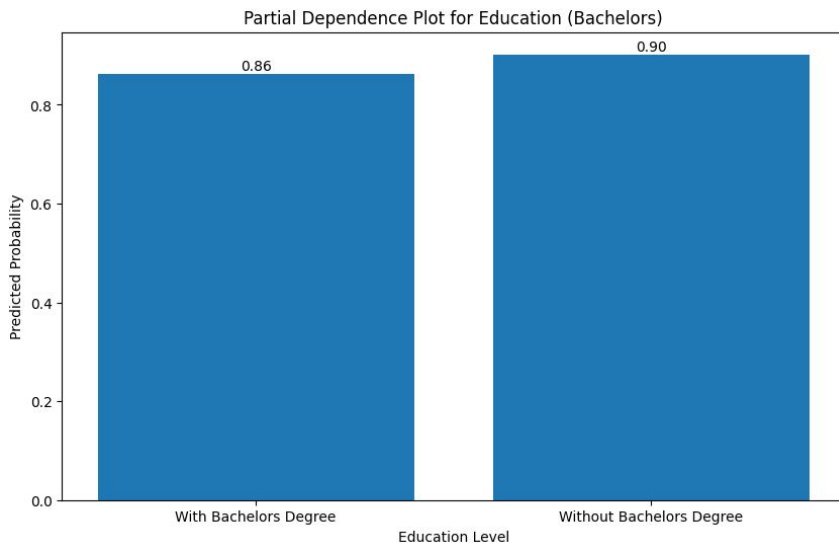
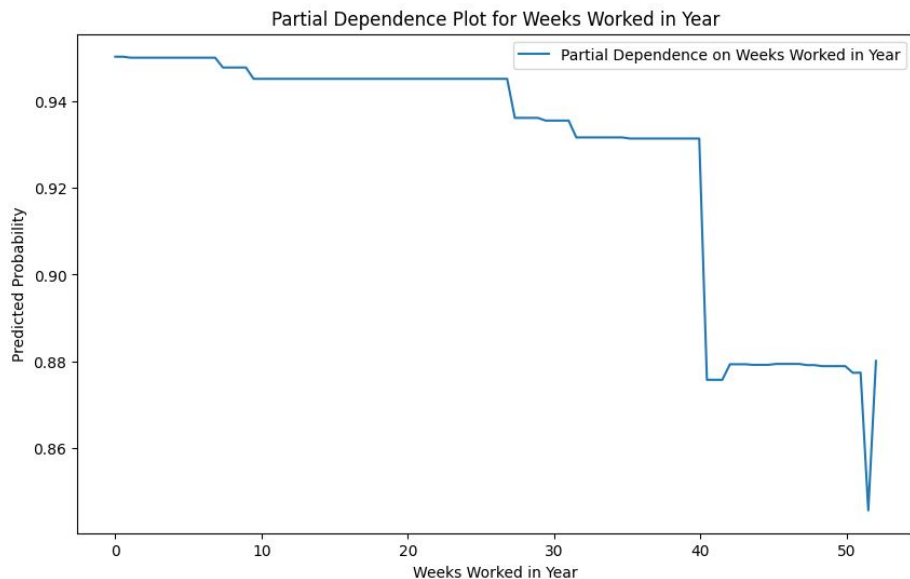


Employees at larger employers less likely >\$50k

- Likelihood of >\$50k **decreases** as the size of one's employer reaches 450 employees



Weeks Worked a Year and Bachelor's Holders



Recommendations

- Enable / encourage employees to remain employed 52 weeks / year
 - Encourage full-year employment and maximized work week
 - Seasonal employers should creatively find employment for temps
- Promote higher education
 - At least at the bachelor's level, employees have a higher probability of earning >\$50k
- Personalized Development Plans (PDPs)
 - PDPs should be based on employee's current skills and their personal as well as business gaps
- Leverage predictive analytics for major employment decisions
 - Promotions, salary changes, development / coaching opportunities
 - And hiring

Considerations for improvements

- Identify how changes in data entries from '94' to '95' elevated, lowered, or retained one's income bracket
- Family structure opens other possibilities
- Feature engineering
 - Grouping highschool and middle school dropouts
 - Grouping college graduates
 - Scaling variables differently
- More models, more optimization
 - LightGBM
 - KNN - nearest neighbor analysis
 - Deep Learning

Q&A