

Employing Machine Learning Methods in Predicting Incomes of Recent College Graduates

Using a principled machine-learning approach, we predict recent college graduates' earnings using data from the College Scorecard. These predictions are estimated using elastic net regularization and the random forest algorithm, regression-based methods adept at producing parsimonious statistical models and reducing bias. Our results support the predictive capabilities of institutional characteristics like school classification, overall debt repayment rates and family income on recent graduate earnings.

QUESTIONS

1. What institutional/program level variables are most predictive of median first-year earnings?
2. What are the patterns of association (positive/negative) for these variables?

DATA

College Scorecard

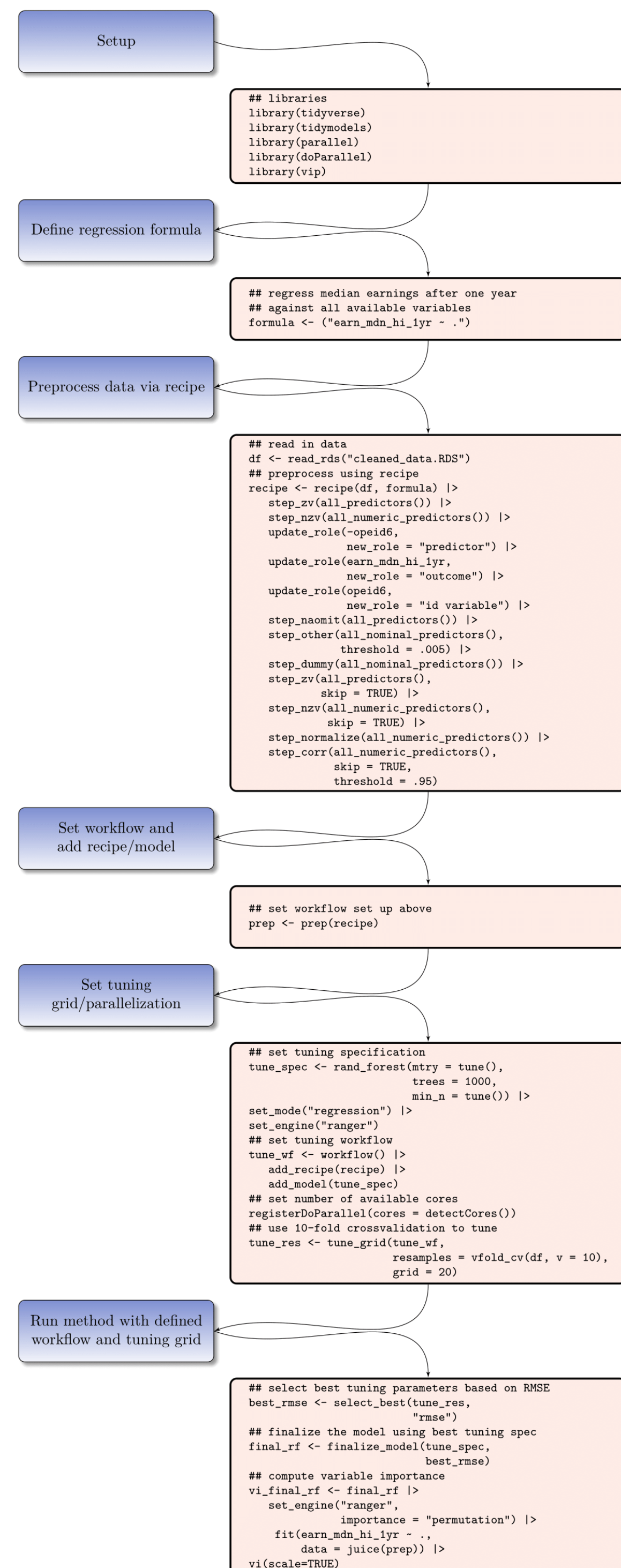
- 2019–2020
- Includes over 2,000 variables for 6,700 accredited higher education institutions (HEIs) in the U.S.
- Significant portions of the data are privacy suppressed

American Community Survey

- 2015–2019
- Geography of interest: County level (using FIPS)
- Matching HEIs with unique county FIPS codes
- Helps recover information lost from missing data in the Scorecard

WORKFLOW

Using a Tidymodels workflow (Kuhn and Silge, 2022) in the R language, we preprocess our data and perform analyses following a standardized recipe structure that allows for transparent parameter tuning and model adjustments. Below is example code for our random forest model.



MODELS

Elastic net regularization (Zou & Hastie, 2005)

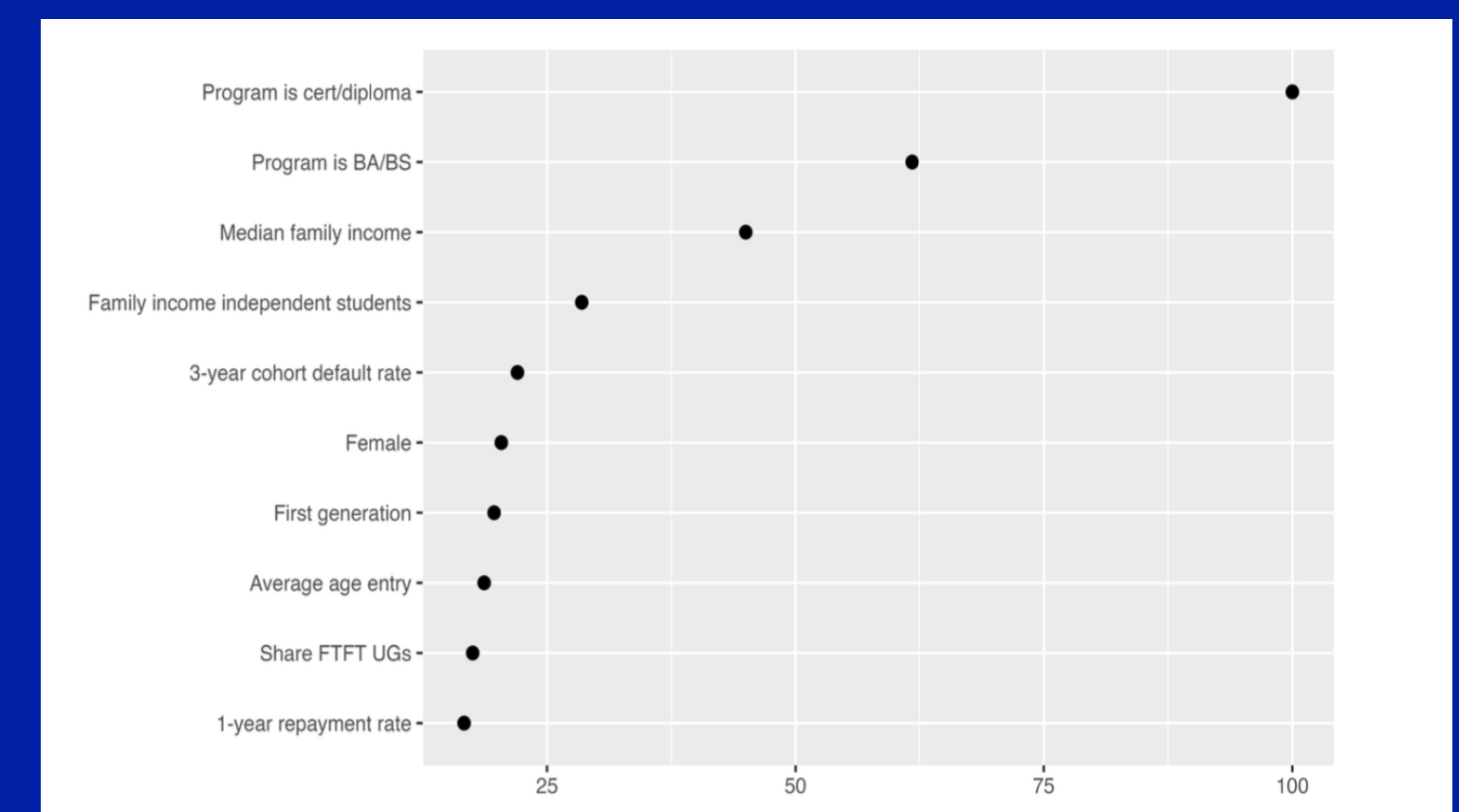
- Combines LASSO/ridge regression penalties to mitigate issues of bias and variance introduction in your statistical model

Random forest

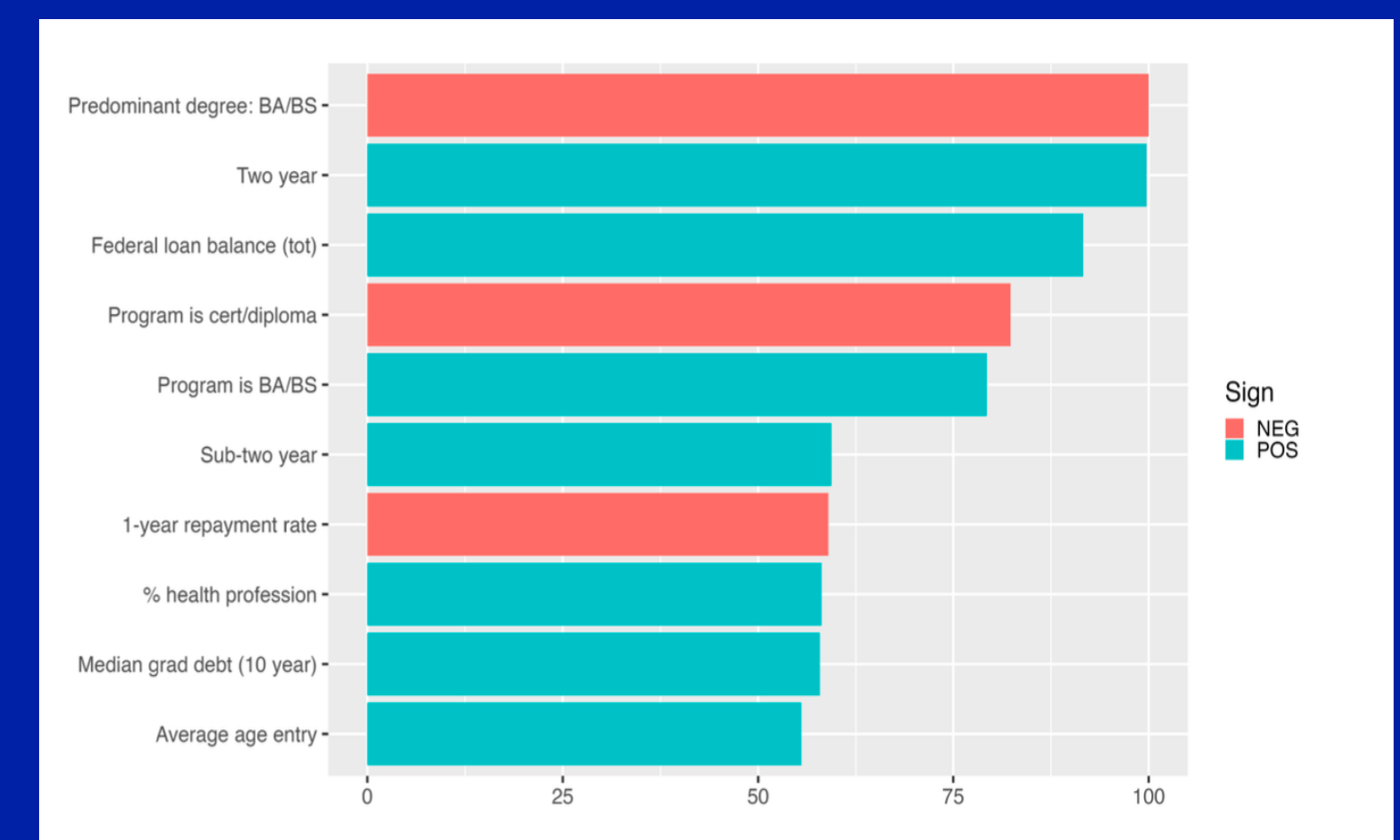
- Randomly samples covariates in a regression model, splits observations on covariates to maximize reduction of RMSE

RESULTS

RQ1 (Variable Importance)



RQ2 (Elastic Net Model)



AUTHORS

Olivia Morales, University of Florida
Benjamin Skinner, University of Florida
William Doyle, Vanderbilt University

