Georgia Institute of Technology

Introduction

Goal: Identify and model relationships race, gender, title, and rating have with faculty salaries at research universities

Reason: Salaries vary a lot past different factors, which factors are most determining?

Data Aggregation

Base Data: 7-8 fields for 497,834 observations, public institution data from TN, IL, CA, GA, NC

Macro Data: Analysis across all universities Micro Data: Analysis across Georgia Tech

Race + Gender: Uses demographic distribution to best guess off last and first name respectively (ethnicolr, genderguesser)

Rating: Scraped from headless browser into CIOS/smartevals. Systematic gathering of by professor evaluation data



Endpoint URL

Determining Factors and Patterns in University Faculty Salary Differences with Regressive and Correlative Analysis Tyler Branscombe & Edmund Chen

Data Processing

Role Normalization: Fuzzy matching and generalization of professor roles into four main roles (Assistant, Associate, Adjunct, Full) • ie. VST ADJ PROF \rightarrow Adjunct

Rating Normalization: Rating data many dimensions, normalized to one overall effectiveness course score, further reduced out of 5 and matched to Georgia Tech professors (roughly ~70% matches)



Research Impact: Challenges matching and gathering enough significant datapoint coverage

Regression Modeling

Modeling: (for 3 subsets of our data) Linear, Lasso, Ridge, Neural Network, Logistic $Y_i = f(X_i, \beta) + e_i$

Evaluation of models: R-Squared, RMSE, MAE

Evaluation of Variables: P-value of linear model for signifigance Coefficient of Linear model to see impact

* Method	R2 [‡]	RMSE
linear model	0.393053478600865	0.02999788
Neural Network	0.404417784088155	0.02971537
Ridge Regression	0.393132238841919	0.02999458
Lasso Regression	0.385788195522131	0.03023150
Logistic Regression	0.393015767253654	1.13637920



