# Predicting Earnings after Graduation

Jenna Kutz, Percy Musiiwa, Matt Jackson

March 12, 2024

#### Problem Statement

Education is the cornerstone of a successful career. In a fast-moving society, the vast majority of potential employers seek to hire individuals with postsecondary degrees and certificates.

According to the Association of Public & Land-Grant Universities,

"The evidence that a college degree significantly improves one's employment prospects and earnings potential is overwhelming... On an annual basis, median earnings for bachelor's degree holders are \$36,000 or 84 percent higher than those whose highest degree is a high school diploma."

As a cohort that is soon to complete a six-month certificate, we (at least intuitively) understand the value of education. But how can a potential student be more deliberate if the goal is to maximize income over the course of a career?

### Project Overview

#### Goals:

- Provide insights on post-graduate earnings from the <u>data.gov college scorecard</u>
- Create a multiple linear regression model that will provide insight into the universities and features that correlate to higher earnings
- Utilize tools from previous modules to complete the analysis:
  - o PySpark / SQL
  - O Data cleansing, standardization tools (Standard Scaler)
  - Linear regression libraries (sklearn)
  - O General data analysis tools (Python, Pandas, etc.)

## Data Exploration and Familiarization

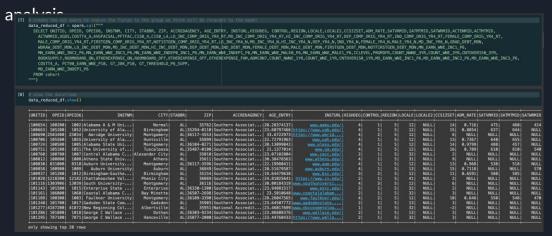
- Dataset contained ~6500 rows and ~2900 columns
- A data dictionary was included (very helpful!)
- Data categorized into broad categories:
  - Academics
  - Admissions
  - o Aid
  - Completion
  - o Cost
  - o Earnings
  - O Repayment
  - o School
  - Student

For our analysis, we tried to build as much familiarity with the cost, earnings, school and student domains.

### Data Prep and Ingestion

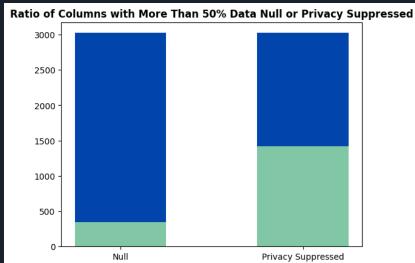
The data.gov College Scorecard was broken into multiple files, each ~250MB to prep the data for ingestion we took these steps:

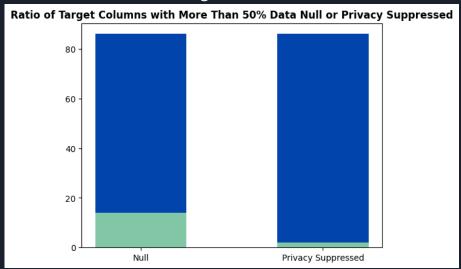
- Created a zip file (~24MB) of the most recent cohort that can be loaded to our GitHub repo
- Aligned on the use of Google Colab as our notebook environment
- Aligned on using Spark to ingest and perform the initial rendering of the data
- Collectively selected ~80 fields that would serve as our starting point for feature



### Data Cleanup & Feature Selection

At this point, we decided that we would focus on creating a multiple linear regression model. To prep the data for this type of analysis, we'd need to remove string and null values from the





#### Target variable:

Median earnings of independent students working and not enrolled 6 years after entry

# Model Training and Optimization

- After converting the Spark dataframe to a Pandas dataframe, we ran the .corr() method on the remaining features to gather a baseline for relatedness to our target variable.
- After importing the sklearn modules, we ran the initial data.
- At first, our model was only performing in the 55-65% range. We had to find more feature variables that were better correlated to the target.
- Optimizing the model took approximately 10-15 attempts but we were able to get it above 80%

#### Correlation to Target

```
model_fields.corr()['MD_EARN_WNE_INDEP1_P6']
HIGHDEG
                               0.524405
AGE ENTRY
                              -0.231888
NOTFIRSTGEN_DEBT_MDN
                              0.559862
FIRSTGEN DEBT MDN
                               0.571523
MALE_DEBT_MDN
                               0.601051
FEMALE DEBT MDN
                               0.549165
IND DEBT MDN
                               0.500618
DEP DEBT MDN
                               0.592892
HI INC DEBT MDN
                               0.557926
MD INC DEBT MDN
                              0.552280
LO INC DEBT MDN
                               0.544745
WDRAW_DEBT_MDN
                               0.421478
GRAD DEBT MDN
                               0.376582
MD_INC_YR6_N
                               0.002820
MALE YR4 N
                               0.007979
FEMALE YR4 N
                              -0.010155
IND YR4 N
                              -0.071919
DEP_YR4_N
                               0.094483
HI INC YR4 N
                               0.301812
MD INC YR4 N
                              0.057178
LO_INC_YR4_N
                              -0.150689
NOT1STGEN COMP ORIG YR4 RT
                               0.595411
FIRSTGEN_COMP_ORIG_YR4_RT
                               0.624209
MALE COMP ORIG YR4 RT
                               0.621164
FEMALE COMP ORIG YR4 RT
                               0.620534
IND COMP ORIG YR4 RT
                               0.622742
DEP_COMP_ORIG_YR4_RT
                               0.589223
HI INC COMP ORIG YR4 RT
                               0.503009
MD INC COMP ORIG YR4 RT
                               0.544615
LO INC COMP ORIG YR4 RT
                               0.640177
COUNT NWNE 1YR
                              0.036239
COUNT_NWNE_3YR
                              -0.058764
CCSIZSET
                               0.512436
UGDS
                               0.152350
COSTT4 A
                               0.549476
PCT90_EARN_WNE_P10
                               0.733801
GT 28K P10
                               0.803089
GT THRESHOLD P6 SUPP
                               0.874590
MD EARN WNE INDEP1 P6
                               1.000000
Name: MD EARN WNE INDEP1 P6, dtype: float64
```

#### Output and Conclusions

```
score = model.score(X, y, sample_weight=None)
mse = mean_squared_error(y_test, predictions)
rmse = np.sqrt(mse)
r2 = r2_score(y_test, predictions)
print(f"The score is {score}.")
print(f"The mean squared error is {mse}.")
print(f"The root mean squared error is {rmse}.")
print(f"The r2 score is {r2}.")
# print(f"The standard deviation is {std}.")
The score is -719897458.3919673.
The mean squared error is 24271349.987317916.
The root mean squared error is 4926.5961867518545.
The r2 score is 0.8330727452894079.
```