## Project Overview

# I. Definition

## Project Overview

Predicting student future earnings based on college.

Department of Education collects data from universities on students who apply for financial aid, this is publicly available data that can help students determine best college selection. Choosing a college to attend can be difficult task given vast options available. This choice is ever more important given that the cost of tuition in private colleges is over 50,000 per year.

In this paper explore how various college characteristics influence future student earnings and attempt to answer if students attending private colleges have an advantage in future earning over public colleges, or is the situation reversed given that students in public colleges have much lower student loans.

This project analyzes data from [US department of education](#), which collets various data points on schools, admissions and tuition costs, student demographics financing and post education earnings. I will create a model capable of predicting future student's earning based on the available data.

## Problem Statement

The goal is to analyze college score-card data provided through US Department of Education. Data is provided for period of 1996-2014. Data is provided in large number of tab delimited files along with column mapping. To enable easier filtering and exploration of data I need to:
1. Download and preprocess the DOE Score Card Text data. Import data into SQL Database
2. Using classifier explore data to determine important features
3. Visualize important features to verify that there is a relationship with response data
3. Train different classifiers measuring predictive accuracy on Test data
4. Evaluate results and further refine features

Hypotheses

My hypothesis is that we can predict future income 6 years after graduation based on input data namely SAT score of the student, Type of Degree, etc.

Solution

To solve this problem of predicting future student income I am going to use supervised learning   regression classifiers. Classifier will be able to predict future earnings given input data. To pick best algorithm I am going to be using Mean Square Error and $R^2$ Accuracy score.

# Metrics

Due to the large number of features available in the dataset (1731 columns, out of this 397 are numerical) we need a way either programmatically or heuristically picking important variables. Because a programmatic approach is likely going to pick features which are most correlative to the response variable I choose to use the combination of the two approaches. Start by picking variables that related research has shown as important factors to student success, and use programmatic evaluation to confirm importance.

To reduce the feature set to more manageable data set of about 10%

columns I will use two RFE and Feature importance.

- Recursive feature elimination (RFE) for feature selection which selects features by recursively considering smaller and smaller sets of features. RFE outputs list of features according to their predictive importance (of future earnings)
- Lasso as a regression method penalizes extra variables producing list of coefficients for each variable. Higher the coefficients, higher is the importance of the variable
- Feature Importance. Methods that use ensembles of decision trees (like Random Forest or Extra Trees) can also compute the relative importance of each attribute. These importance values can be used to inform a feature selection process.

Predictive Model Evaluation:
- Mean squared error, measuring distance or difference of prediction of data the model has not seen (test data) and true test data. This metric is used because it penalizes large differences between predicted and true values in both positive and negative.   Due to the size of the MSE being in Millions I will be using root of MSE, or RMSE.
- $R^2$ Accuracy score measuring (**1 - u/v**), where **u** is the regression sum of squares ((y_true - y_pred) ** 2).sum() and v is the residual sum of squares ((y_true - y_true.mean()) ** 2).sum()

# II. Analysis
*(approx. 2-4 pages)*
## Data Exploration
Given that the data is very wide (1731 columns) its immediately apparent that in order to have a manageable analysis I need to reduce this dataset. Based on heuristics on which information is likely to influence future earning and Looking at the previous research on the subject such as Brookings Institute.   I choose to concentrate on following part of the data set

| md_earn_wne_p6 | Median earnings of students working and not enrolled 6 years after entry |
|---|---|

| | |
|---|---|
| PREDDEG | Predominant undergraduate degree awarded<br>0 Not classified<br>1 Predominantly certificate-degree granting<br>2 Predominantly associate's-degree granting<br>3 Predominantly bachelor's-degree granting<br>4 Entirely graduate-degree granting |
| HIGHDEG | Highest degree awarded<br>0 Non-degree-granting<br>1 Certificate degree<br>2 Associate degree<br>3 Bachelor's degree<br>4 Graduate degree |
| CONTROL | Control of institution (Private NP, Private P, Public) |
| NUMBRANCH | Number of branch campuses |
| AVGFACSAL | Average faculty salary |
| ADM_RATE | Admission rate |
| SAT_AVG | Average SAT equivalent score of students admitted |
| TUITFTE | Net tuition revenue per full-time equivalent student |
| UGDS | Enrollment of undergraduate certificate/degree-seeking students |
| UGDS_NRA | Total share of enrollment of undergraduate degree-seeking students who are non-resident aliens |
| PPTUG_EF | Share of undergraduate, degree-/certificate-seeking students who are part-time |
| UG25abv | Percentage of undergraduates aged 25 and above |
| PAR_ED_PCT_1STGEN | Percentage first-generation students |
| DEP_INC_AVG | Average family income of dependent students in real 2015 dollars. |
| IND_INC_AVG | Average family income of independent students in real 2015 dollars. |
| COMP_ORIG_YR2_RT | Percent completed within 2 years at original institution |
| WDRAW_ORIG_YR2_RT | Percent withdrawn from original institution within 2 years |
| ENRL_ORIG_YR2_RT | Percent still enrolled at original institution within 2 years |
| COMP_ORIG_YR4_RT | Percent completed within 4 years at original institution |
| WDRAW_ORIG_YR4_RT | Percent withdrawn from original institution within 4 years |
| ENRL_ORIG_YR4_RT | Percent still enrolled at original institution within 4 years |
| OVERALL_YR2_N | Number of students in overall 2-year completion cohort |
| OVERALL_YR3_N | Number of students in overall 3-year completion cohort |
| OVERALL_YR4_N | Number of students in overall 4-year completion cohort |
| OVERALL_YR6_N | Number of students in overall 6-year completion cohort |
| OVERALL_YR8_N | Number of students in overall 8-year completion cohort |
| count_nwne_p6 | Number of students not working and not enrolled 6 years after entry |
| DEBT_MDN | Median debt, suppressed for n=30 |
| GRAD_DEBT_MDN | Median debt of completers, suppressed for n=30 |

I called this subset of dataset data_reduced having 28281 rows and 31 columns out of which 3 were categorical and 28 numerical.

Number of columns have PrivacySuppressed instead of values, In order for linear models to work correctly all of the data has to be of the same type (numeric) in a single column. I replaced PrivacySuppressed with NaN (Null).

Another requirement of linear models is that it requires non Null data. I originally chose to fill the null values with 0 which had similar results as mean. However as pointed out in refinement section this led to significantly changes in the models to the point where prediction logic produced didn't make sense. Instead I chose to drop all the null values. Doing this further reduces the dataset to 3,813 rows and 31 columns out of which 3 were categorical and 28 numerical.

Data Statistics:

| | md_earn_wne_p6 | NUMBRANCH | AVGFACSAL | ADM_RATE | SAT_AVG | TUITFTE |
|---|---|---|---|---|---|---|
| count | 3,813 | 3,813 | 3,813 | 3,813 | 3,813 | 3,813 |
| mean | 35,929.16 | 1.36 | 6,809.88 | 0.67 | 1,043.97 | 9,465.49 |
| std | 8,127.25 | 1.95 | 1,684.15 | 0.18 | 114.15 | 5,682.89 |
| min | 15,600 | 1 | 2,794 | 0.06 | 725 | 359 |
| 0.25 | 31,100 | 1 | 5,652 | 0.57 | 971 | 4,900 |
| 0.50 | 35,100 | 1 | 6,499 | 0.69 | 1,030 | 8,359 |
| 0.75 | 39,500 | 1 | 7,713 | 0.79 | 1,104 | 12,749 |
| max | 118,300 | 23 | 15,922 | 1 | 1,491 | 46,776 |

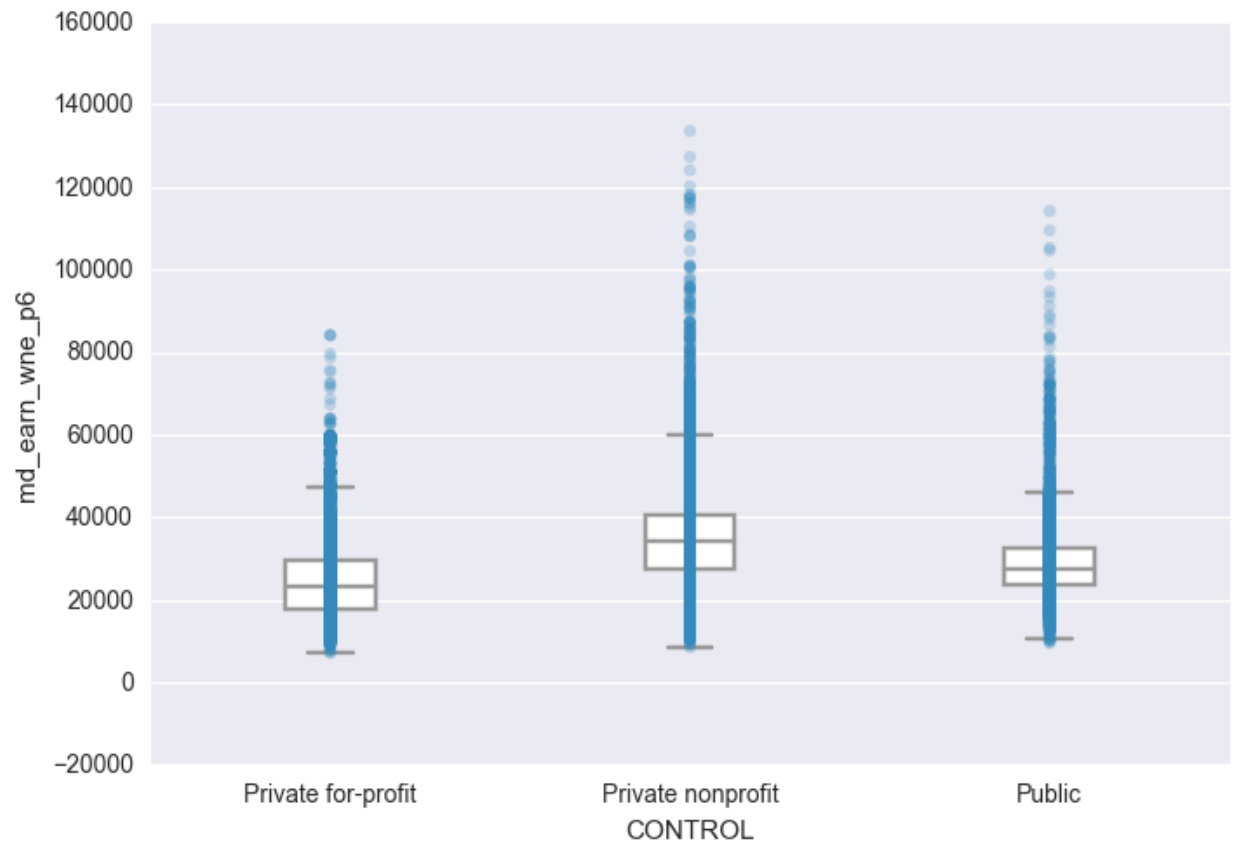| | IND_INC_AVG | COMP_ORIG_YR2_RT | WDRAW_ORIG_YR2_RT | ENRL_ORIG_YR2_RT | CO |
|---|---|---|---|---|---|
| count | 3,813 | 3,813 | 3,813 | 3,813 | 3, |
| mean | 28,764.58 | 0.11 | 0.18 | 0.44 | 0. |
| std | 9,728.64 | 0.07 | 0.08 | 0.11 | 0. |
| min | 904.36 | - | 0.01 | - | 0. |
| 0.25 | 22,297.59 | 0.07 | 0.12 | 0.38 | 0. |
| 0.50 | 27,012.58 | 0.10 | 0.17 | 0.44 | 0. |
| 0.75 | 34,323.16 | 0.14 | 0.22 | 0.51 | 0. |
| max | 73,242.66 | 0.66 | 0.59 | 0.79 | 0. |

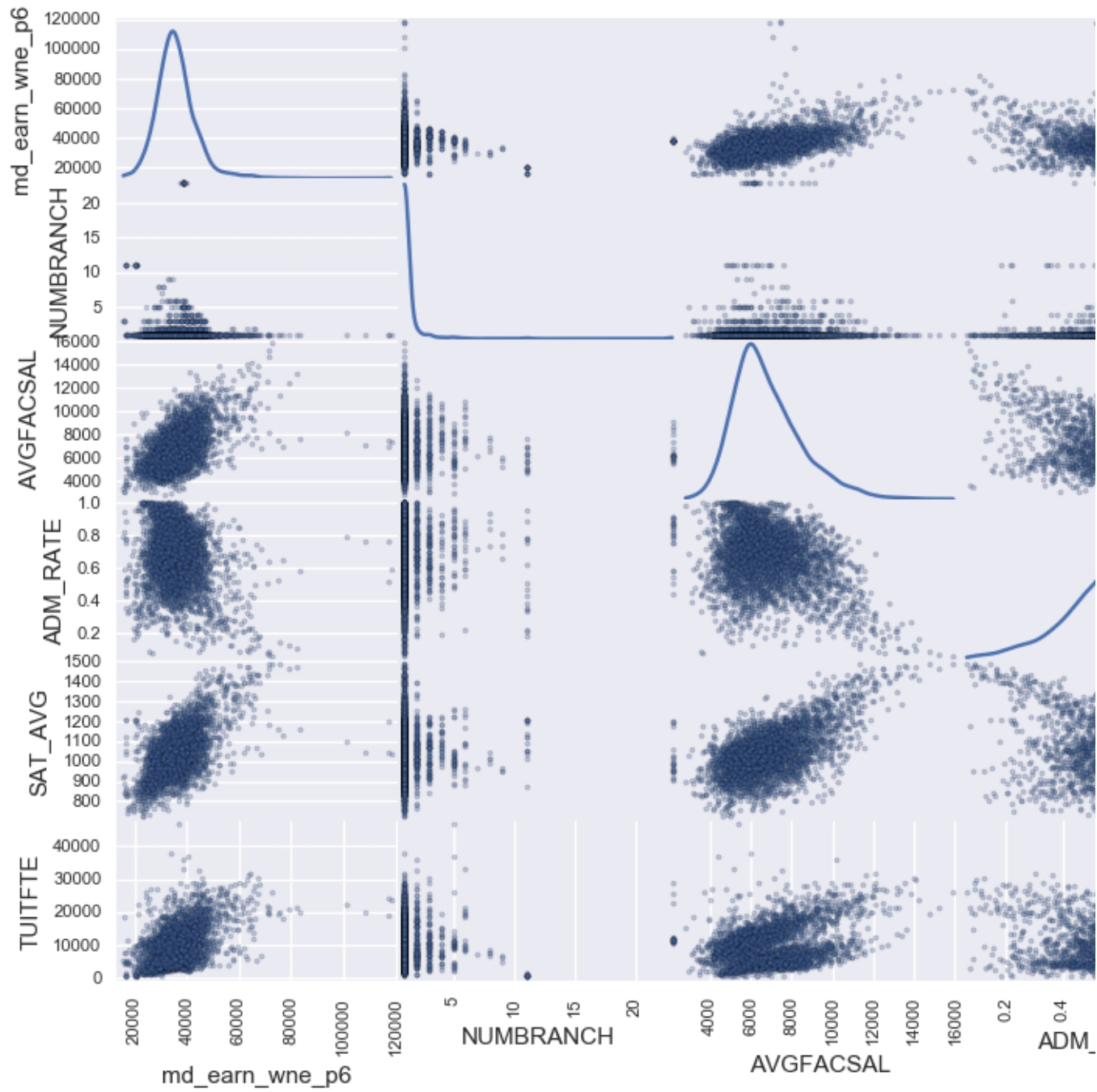| | OVERALL_YR4_N | OVERALL_YR6_N | OVERALL_YR8_N | count_nwne_p6 | DEBT_MDN | G |
|---|---|---|---|---|---|---|
| count | 3,813 | 3,813 | 3,813 | 3,813 | 3,813 | 3 |
| mean | 1,653.48 | 1,547.43 | 1,431.64 | 189.71 | 13,923.51 | 1 |
| std | 1,856.01 | 1,719.71 | 1,609.41 | 477.19 | 3,577.14 | 4 |
| min | 75 | 52 | 18 | 3 | 2,624 | 5 |
| 0.25 | 546 | 514 | 472 | 50 | 11,625 | 1 |
| 0.50 | 998 | 949 | 877 | 96 | 14,000 | 1 |
| 0.75 | 2,044 | 1,935 | 1,780 | 203 | 16,250 | 2 |
| max | 20,854 | 14,293 | 13,852 | 8,098 | 25,750 | 3 |

Observation from statistics:

- Salary after 6 years of graduation is normally distributed with long tail. Median salary is 34K and Maximum salary is 118K. 75% of graduates earn 39.5K per year.
- Average faculty salary is also normally distributed with 75% of the faculty receiving 7,713 .
- Admittance average rate is 67% , which is slightly higher than expected, with very few schools having 100% admittance rate.
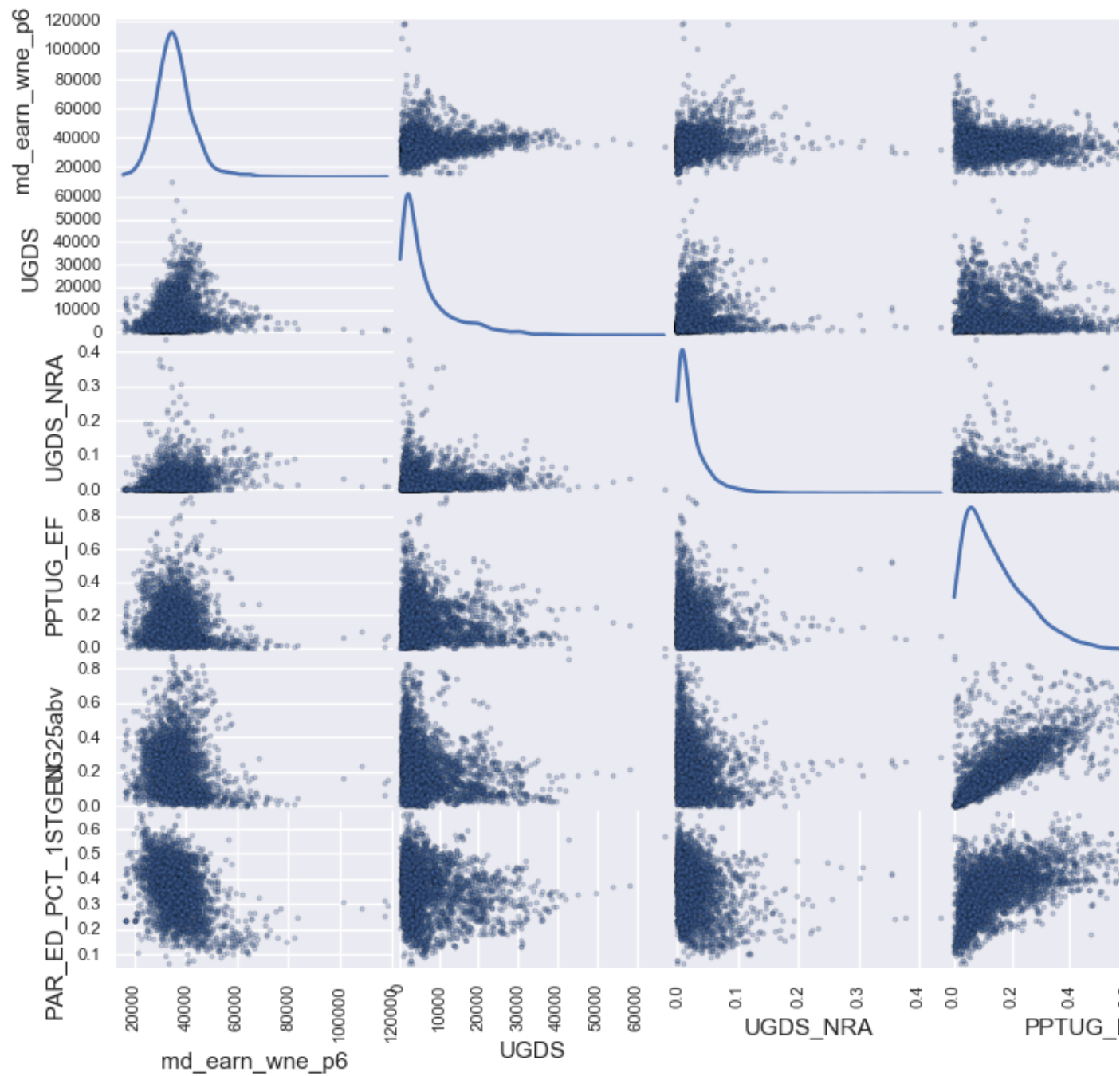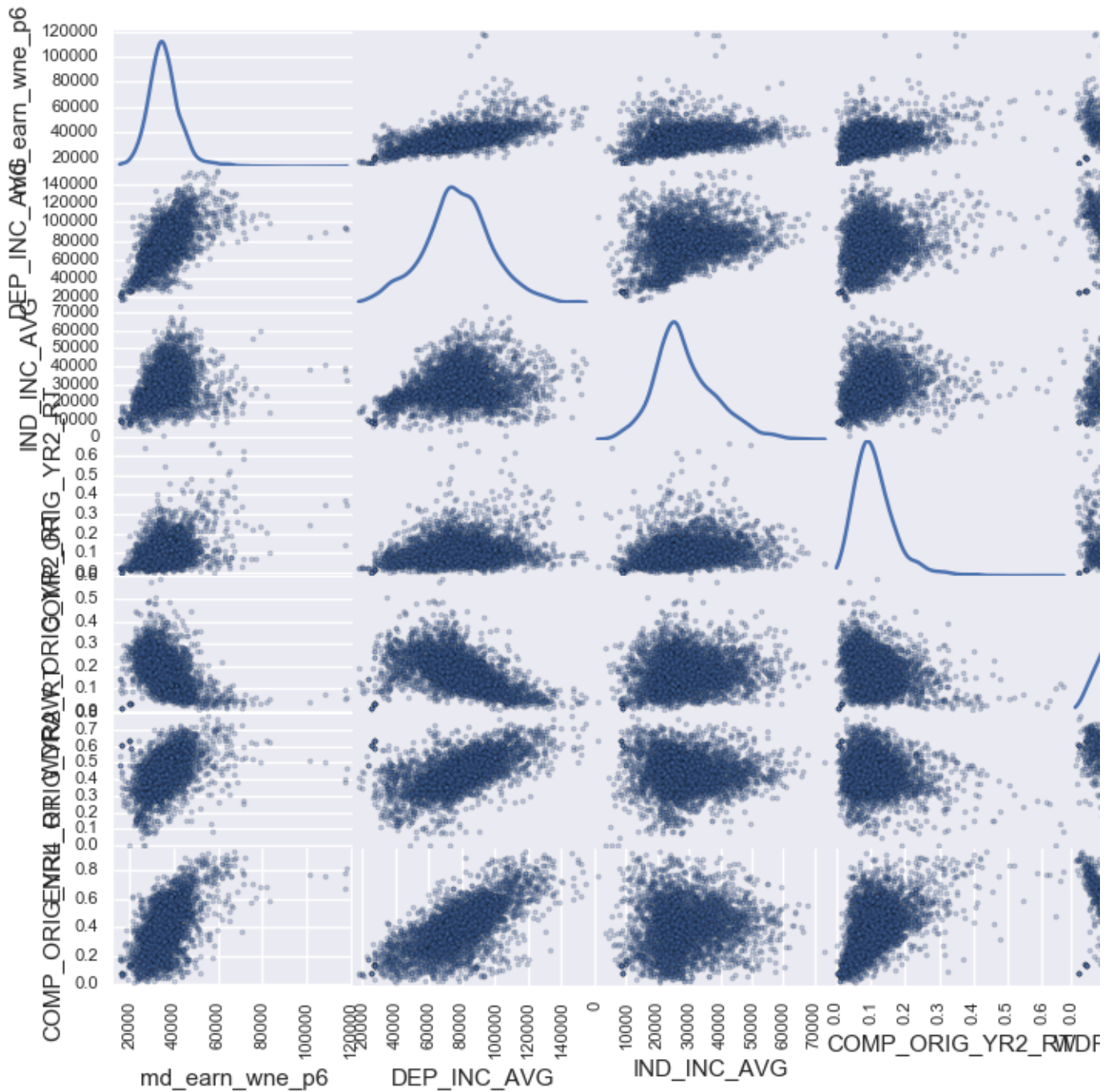- UGDS and

# Exploratory Visualization
Response Variable

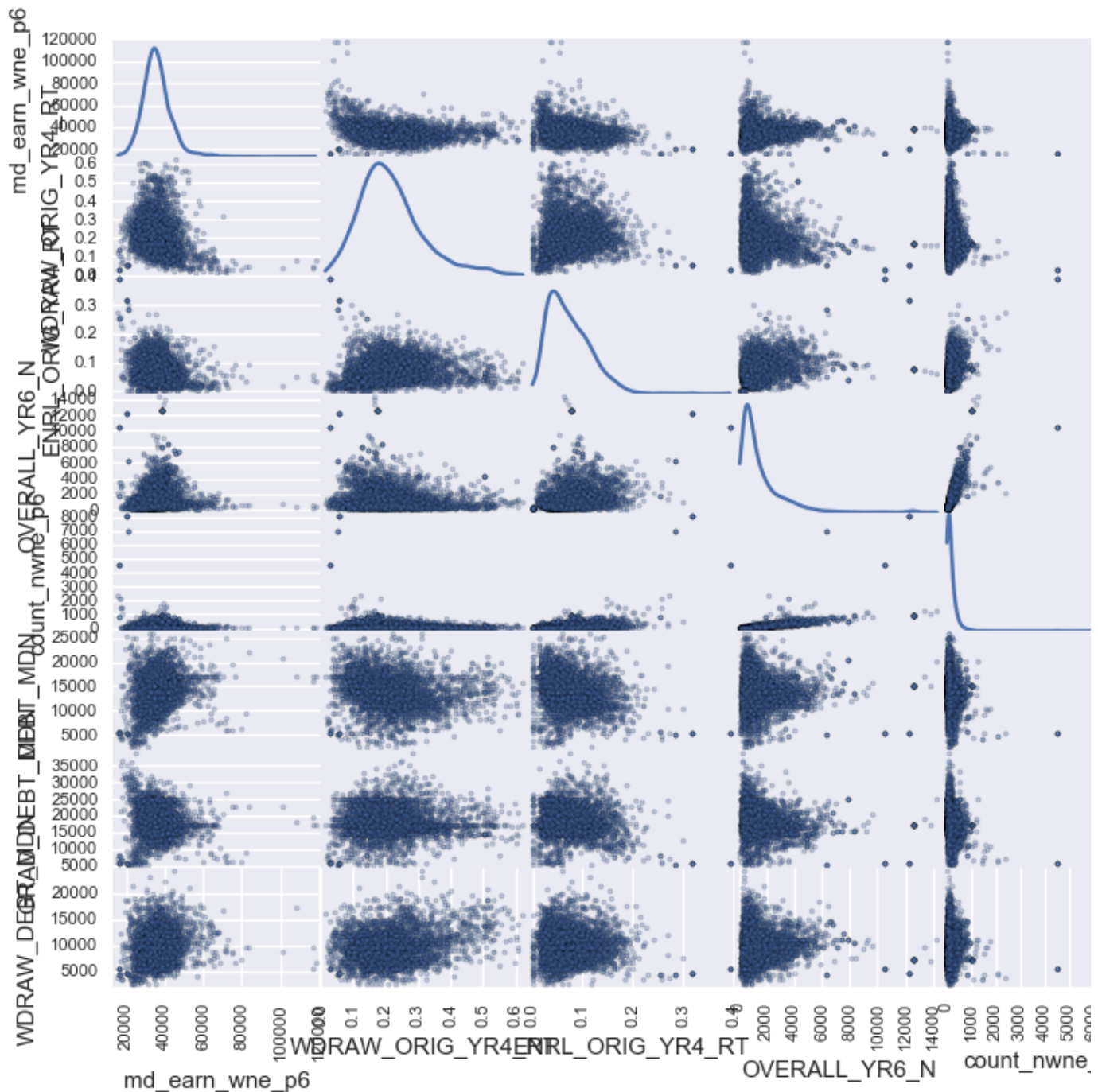- 'md_earn_wne_p6', 'NUMBRANCH', 'AVGFACSAL', 'ADM_RATE', 'SAT_AVG', 'TUITFTE

- 'md_earn_wne_p6', 'UGDS','UGDS_NRA', 'PPTUG_EF', 'UG25abv',
  'PAR_ED_PCT_1STGEN'

- 'md_earn_wne_p6', 'DEP_INC_AVG','IND_INC_AVG', 'COMP_ORIG_YR2_RT', 'WDRAW_ORIG_YR2_RT', 'ENRL_ORIG_YR2_RT', 'COMP_ORIG_YR4_RT'

- 'md_earn_wne_p6', 'WDRAW_ORIG_YR4_RT', 'ENRL_ORIG_YR4_RT', 'OVERALL_YR6_N', 'count_nwne_p6', 'DEBT_MDN', 'GRAD_DEBT_MDN', 'WDRAW_DEBT_MDN'

- 

Graphs above show comparison of each variable to our response variable.
Kernel Distribution is shown Diagonally.
I chose to compare the response variable to each feature which shows
collinearity between them Strongest relationship can be observed between
future income and

- PAR_ED_PCT_1STGEN      Percentage first-generation students
- AVGFACSAL      Average faculty salary

- SAT_AVG   Average SAT equivalent score of students admitted
- DEP_INC_AVG   Average family income of dependent students in real 2015 dollars.
- IND_INC_AVG    Average family income of independent students in real 2015 dollars.

# Algorithms and Techniques

To solve this problem of predicting future student income I am going to use supervised learning   regression classifiers. Classifier will be able to predict future earnings given input data. To pick best algorithm I am going to be using Mean Square Error and Accuracy score.
I will use following regression algorithms

- LinearRegression
- Ridge Regression including Cross-Validation
- Lasso including Lasso Lars Cross-Validation
- RandomForestRegressor

Because Ridge Regression and  Lasso require passing of **alpha parameter**  I will use Cross-Validation to find a best parameter.

**Linear Regression** is useful predictor, but it can suffers from collinearity. With each model that exposes Coefficient estimates I will report on them. Ordinary Least Squares rely on the independence of the model terms. When terms are correlated and the columns of the design matrix X have an approximate linear dependence, the design matrix becomes close to singular and as a result, the least-squares estimate becomes highly sensitive to random errors in the observed response, producing a large variance. This situation of multicollinearity can arise, for example, when data are collected without an experimental design.

**Ridge regression** addresses some of the problems of Ordinary Least Squares by imposing a penalty on the size of coefficients. The ridge coefficients minimize a penalized residual sum of squares.

The **Lasso** is a linear model that estimates sparse coefficients. It is useful in some contexts due to its tendency to prefer solutions with fewer parameter values, effectively reducing the number of variables upon which the given solution is dependent.

A **random forest** is a meta estimator that fits a number of classifying decision trees on various sub-samples of the dataset and use averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is always the same as the original input sample size but the samples are drawn with replacement if bootstrap=True

## Benchmark

Outside benchmark is not available for this problem so I will use liner regression trained on Average SAT score and has very high MSE score and 25% r^2 accuracy.  I am using SAT score because of its strong correlation to response variable and because of assumption that higher SAT score will lead to acceptance into highly ranked school, which according to traditionally held belief will result in high earnings after graduation.

After running Liner Regression on available data following benchmark values are obtained:
Mean squared error root: `6761.79`
Variance score: `0.25`

# III. Methodology
*(approx. 3-5 pages)*
## Data Preprocessing

Due to different scales the data is on preprocessing is required for most regression algorithms. I chose to use centering around the mean as method for data scaling. *Based on the **Data Exploration** I detected some outlier is tuition fee. I will be removing tuition fees higher than 1Million.

Further pre-processing was performed on PREDDEG, HIGHDEG and CONTROL features by `Encode labels with value between 0 and n_classes-1. Original Features were dropped after visualization was performed.`

`Once dataset only contained numerical values, each  feature was filtered for PrivacySuppressed values and where found`

```
PrivacySuppressed was replaced with null.

One of the requirement of regression models is that data
can't contain null values I have dropped nulls.
```

## Implementation

To predict future earnings I chose to use following classifiers

- LinearRegression
- Ridge Regression including Cross-Validation
- Lasso including Lasso Lars Cross-Validation
- RandomForestRegressor

To prevent overfitting and to test accuracy of the models I have separated data into training and test datasets, retaining 40% of data for test. `16286 data points across columns were available for training and 10858 data points for test.`

After running each of the models above following are results of the tests:

| Model | Mean squared error Root | Accuracy |
|---|---|---|
| Benchmark | 6761.79 | 0.25 |
| Linear Regression | 5025.61 | 0.60 |
| Ridge Regression | 5033.65 | 0.60 |
| Lasso | 5025.76 | 0.60 |
| RandomForestRegressor | 4273.54 | 0.72 |

As we can see from the results above, each of the algorithms over performed the benchmark model, which in truth was not difficult to do since accuracy of the benchmark model was only 25%.

## Refinement

In hopes of improve accuracy of the models and since number of models used require parameter to be passed specified I will use Cross-Validation with Ridge Regression and Lasso Lars. For Random forest I will use bootstrapping.

| Model | Mean squared error | Accuracy |
|---|---|---|
| Benchmark | 6761.79 | 0.25 |
| Linear Regression | 5025.61 | 0.60 |
| Ridge Regression | 5033.65 | 0.60 |
| **Ridge Regression Cross-Validation a=0.1** | 5026.42 | **0.60** |
| Lasso | 5025.76 | 0.60 |
| **LassoLarsCV a=0.2** | 5019.12 | **0.60** |
| RandomForestRegressor | 4256.50 | 0.72 |
| **RandomForestRegressor with bootstrap** | 4273.54 | **0.72** |

After comparing the results of final model and examining the feature importance it became clear that filling missing values with 0 was significantly affecting the results.

For example examining feature importance, which I cover later in the section it would appear that the most important variable for predicting future salaries is "Number of students in overall 2-year completion cohort" When I examined the graph it showed that 0 was the highest predictor, this of course doesn't make sense. At this point I have gone back and modified pre-processing step to **Drop rows with null values. This was of course** a major setback as I had to re-evaluate all of the points I have made.

Implementing above algorithms was fairly straightforward and no major issues with coding were discovered.
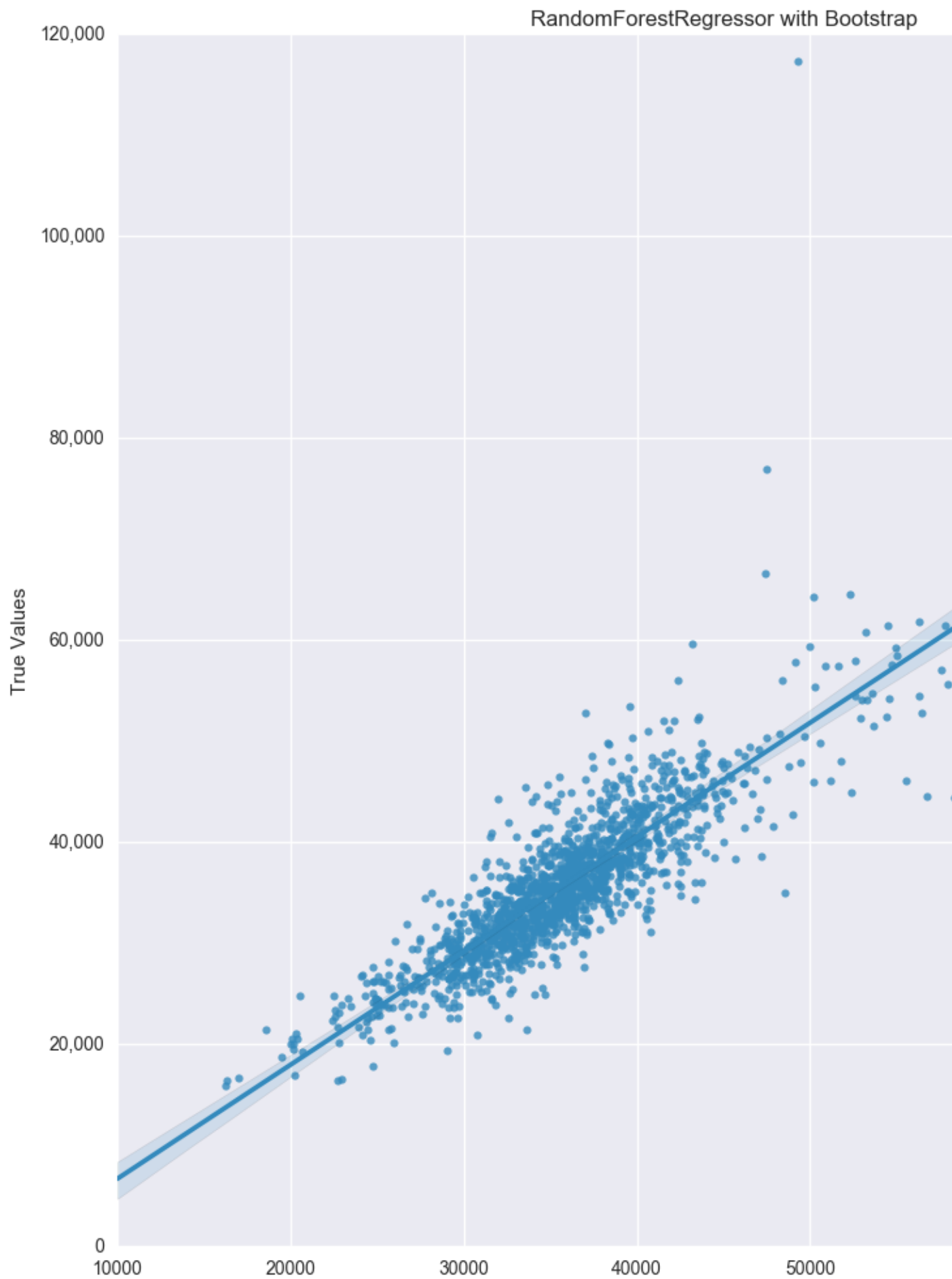
# IV. Results

*(approx. 2-3 pages)*

## Model Evaluation, Validation and Justification

As can be observed from results obtained in model training and refinement **RandomForestRegressor with bootstrap** performs the best in term of accuracy. In visualizing the predicted vs true values we can see that model predicts well for salaries between 25,000 and $50,000 after which data

becomes sparse which results in more inaccurate prediction. The line and the shading represents margin of error of predictions.

I used 100 estimators (The number of trees in the forest). with bootstrap to draw the predictions.  To look at overfitting I looked at out of bag prediction error as s the algorithm builds the forest it remembers the out-of-bag (OOB) prediction error, which is its best guess of the generalization error. Based on OOB error of `0.71` the `RFR is not overfitting.`

# V. Conclusion
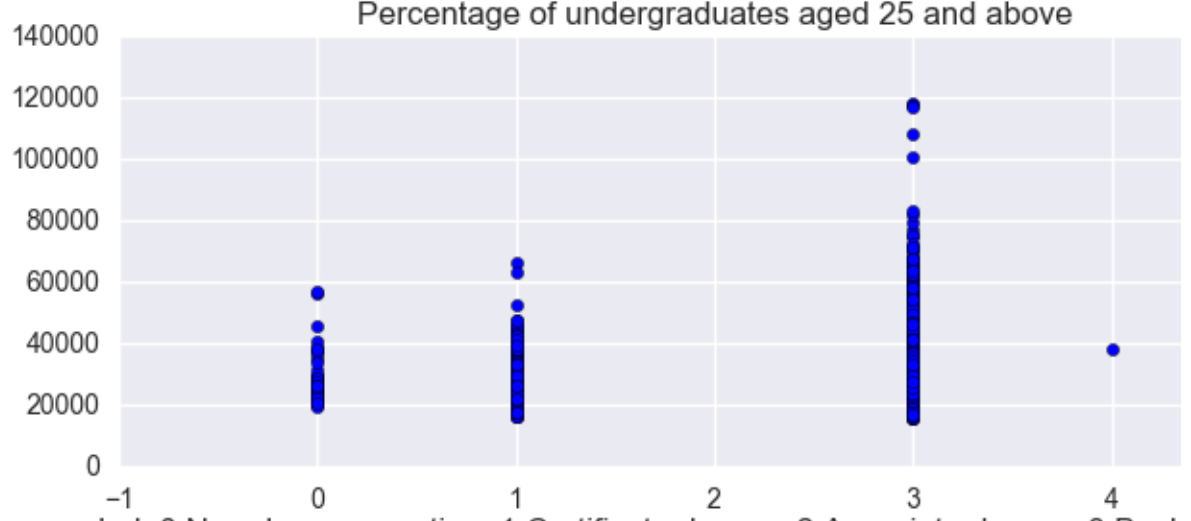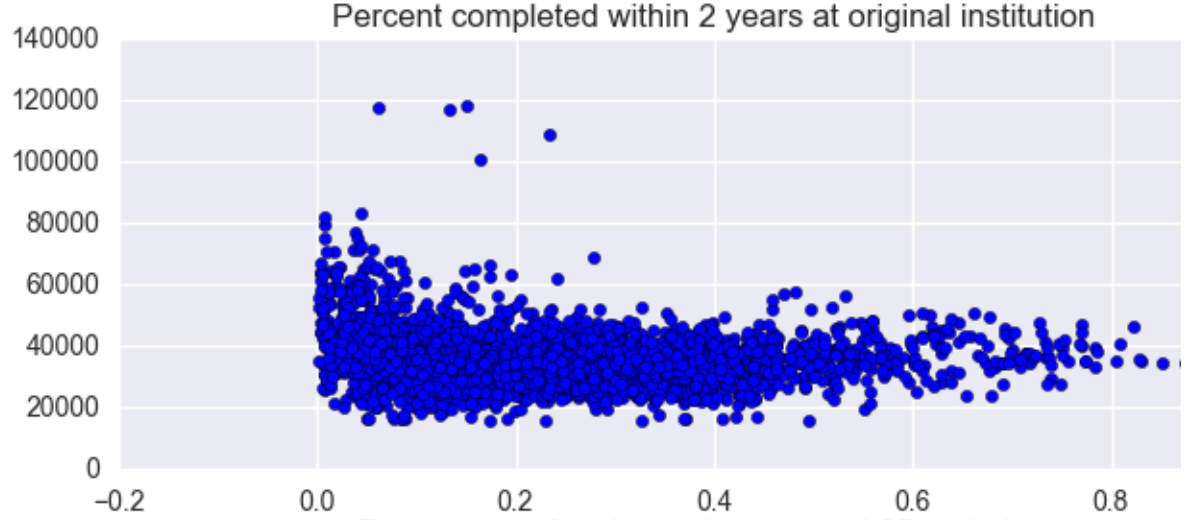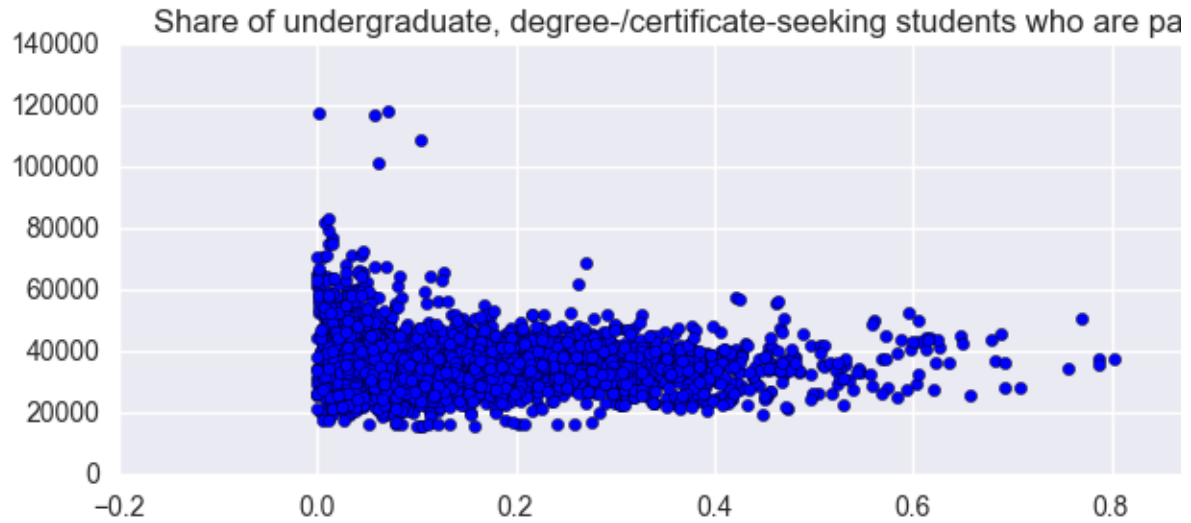
*(approx. 1-2 pages)*

## Free-Form Visualization

Feature Importance
Random Forest Regressors   with boot strapping was shown to have highest accuracy, and while most linear models expose coefficients, Random Forest doesn't but instead it exposes Feature importance. I will examine top 10 features according to their importance.

| Model | RFR |
|---|---|
| PPTUG_EF | 1 |
| COMP_ORIG_YR2_RT | 0.59 |
| UG25abv | 0.25 |
| HIGHDEG_N | 0.25 |

In next series for visualizations I will compare response variable to each of the important features:

Share of undergraduate, degree-/certificate-seeking students who are pa
Percent completed within 2 years at original institution
Percentage of undergraduates aged 25 and above
Highest degree awarded  0 Non-degree-granting  1 Certificate degree  2 Associate degree  3 Bach

What graphs and the important features suggest is that for high post-graduation salary most important features (from the ones examined here are :

- Highest salaries are for
- Share of undergraduate, degree-/certificate-seeking students who are part-time **that are close to 0**. Meaning Students that were not part time were more likely to have higher salaries
- Percent completed within 2 years at original institution **that are close to 0.** Meaning Students who continued studying after 2 were more likely to have higher salaries
- Percent of undergraduates aged 25 and above that at 3 percent or above. Meaning Students who were 25 or older were more likely to have higher salaries
- Highest degree awarded. With Bachelor degree having highest earnings

# Reflection

This was definitely a very challenging project, in many aspects. From data import with large number of files to process, to very large feature space over 1700 columns.

For all of the reasons listed above and in general predicting student income after graduation is a very difficult task, I am glad that I was able to come close to predicting it with 72% accuracy. Looking at the predictions and factors that influence high income after graduation and tying it back to conventional wisdom it makes sense.

What data appears to suggest is concentrate on studying (don't do it part time), Stay in school for more than 2 years and graduate with Bachler degree.

The end to end process was to import the data, filter out missing values, find features to use run through different regression algorithms to find best performing one. Once the best algorithm was picked improving on it by tuning parameters.

# Improvement

There are number of areas of improvements in this project in following

areas:

- Feature Selection. In this project I have analyzed fraction of all the available features. One idea is to programmatically select features to use by starting with an empty model and them looping through list of features adding one by one and based on gain in predictive performance keeping the feature.
- Feature engineering and categorization. Majority of features are non-numerical (397 columns, out of 1731). I would like to transform all the columns to numerical values.
- Missing data. This dataset is very sparse with only 13% of the dataset having all the values for columns I selected. I would like to examine each column and choose appropriate fill value (mean, min, max, 0, etc.)
- Models. There are number of other regression models that I would like to try out. This is in addition to using other models such as Using Neural Networks with Regression that I would like to try out.