

# Student Loan Repayment Prediction

Rojina Deuja, July 2017

## 1. Executive Summary

As a part of the final requirement for the completion of Microsoft Professional Capstone: Data Science, this report is prepared describing the prediction model created for the Capstone project. The datasets required to be explored were provided and the output was thus calculated using the dataset. The data set can be categorized into two different types:

- a. Training set (Values available in train\_values.csv and labels in train\_labels.csv)
- b. Test set values as test\_values.csv
- c. Submission Format with submission\_format.csv.

The major purpose of being provided the dataset is to carry out predictions for Educational Data Mining (EDM) purposes. The dataset is representative of the US student's demographics. It helps to evaluate the student's ability to repay loans that depends on multiple factors. The dataset is obtained from the published information from US Department of Education about student loan. The dataset comprises of a total of 8705 observations or data objects. The test set consists of a total of 6391 unlabeled data that has to be predicted with optimum accuracy, by implementing the most suitable techniques.

## 2. Data Preprocessing

Firstly, the data was explored amongst the 444 features provided in the given data and extracted the most influential attributes for input to the model. To avoid the problems like "the curse of dimensionality" only 83 attributes were chosen as input. Then the data was fit into various machine learning models to choose the one with the best performance in terms of accuracy and low rate of error. Finally, a linear regression model was found to give the most satisfactory output.

The features selected for the model in this analysis were:

Firstly, some essential summary and descriptive characteristics of the data were calculated. Excel was used to quickly compute the following summary statistics for the repayment\_rate:

Mean 47.3709

Median 44.855

Minimum value 5.1627

Maximum value 100.4736

Standard Deviation 20.9876

Some graphs are shown below for insight:

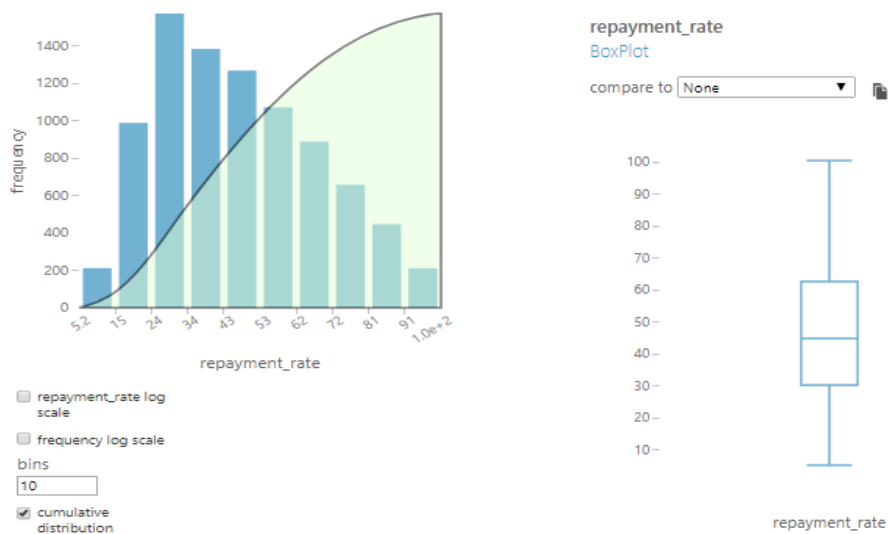


Fig. repayment\_rate column histogram and box-plot

## 2.1 Data Cleaning

One of the major issues in the given dataset was Missing values. Thus, extensive data cleaning was incorporated. The data had two major types of missing values: a. Numeric and b. Categorical variables. The numerical data was cleaned by substituting the missing row with “0”. Similar, the categorical variables were cleaned by substituting the missing row with a value “unknown”.

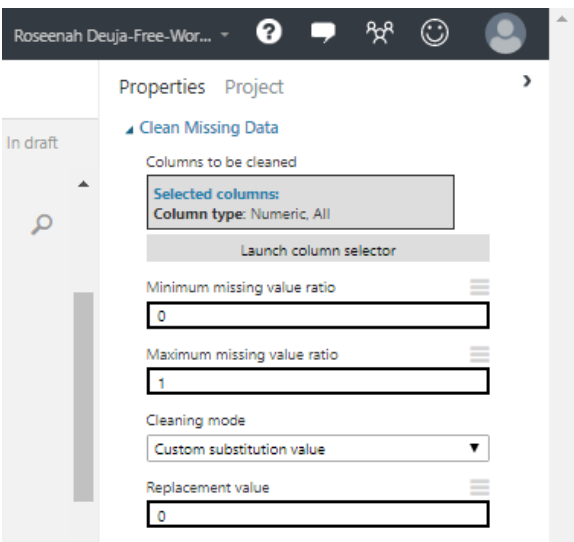


Fig. Handling of missing values (For Numeric)

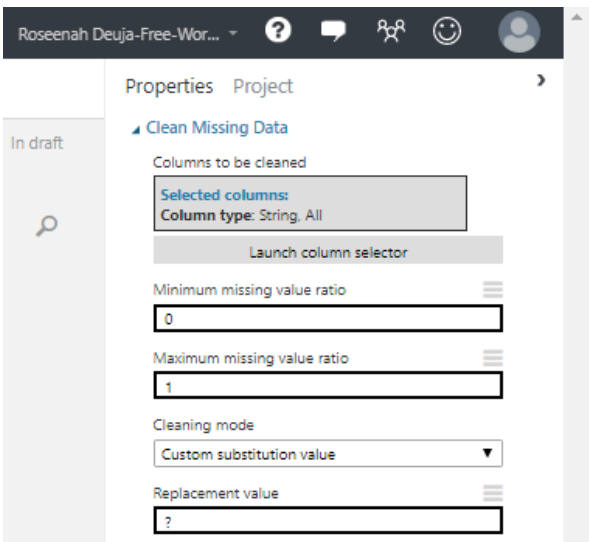


Fig. Handling of missing values (For Categorical)

## 2.2 Data Transformation

Since the different attributes of the data fall in different range of values, it is necessary to transform them to a consistent form. For this Normalization was carried out.

For normalization, Z-Score normalization was used. It is obtained mathematically as:

$$z = \frac{x - \mu}{\sigma}$$

$\mu$  = Mean  
 $\sigma$  = Standard Deviation

Fig. Z-score formula

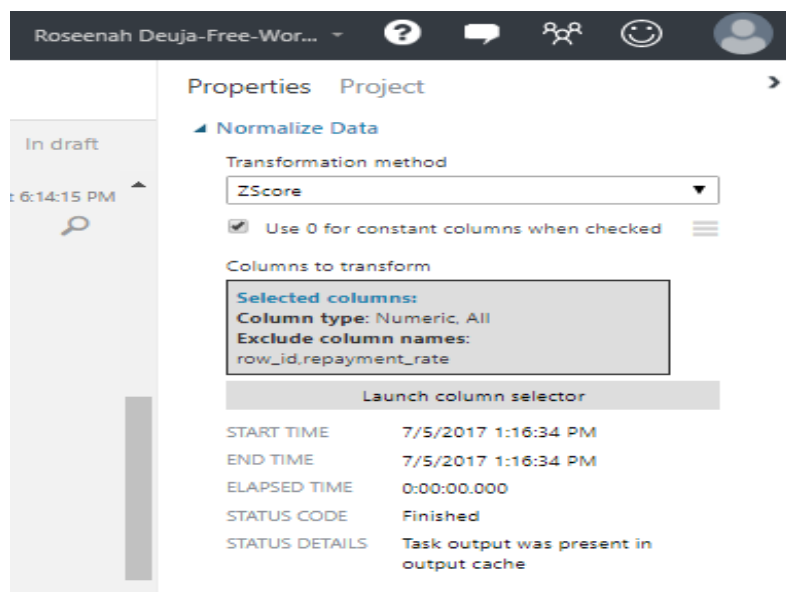


Fig. Normalization by z-score

## 2.3 Data Reduction

The actions preformed under this task are:

- Removing irrelevant attributes: attribute selection, searching the attribute space.
- Principle component analysis (numerical attributes only): After the analysis, it was found that student\_demographics\_avg\_family\_income is the most influential features amongst all features in the dataset. This step helps to determine a lower dimensional space that can best represent the data
- Reducing the number of attribute values: Similarly, the other attributes were also chosen out of the given set. The 83 attributes chosen for analysis are:
  - admissions\_\_act\_scores\_midpoint\_cumulative
  - admissions\_\_act\_scores\_midpoint\_english

3. admissions\_\_act\_scores\_midpoint\_math
4. admissions\_\_act\_scores\_midpoint\_writing
5. admissions\_\_admission\_rate\_overall
6. admissions\_\_sat\_scores\_average\_overall
7. admissions\_\_sat\_scores\_midpoint\_critical\_reading
8. admissions\_\_sat\_scores\_midpoint\_math
9. admissions\_\_sat\_scores\_midpoint\_writing
10. aid\_\_cumulative\_debt\_number
11. aid\_\_federal\_loan\_rate
12. aid\_\_loan\_principal
13. aid\_\_median\_debt\_completers\_overall
14. aid\_\_median\_debt\_dependent\_students
15. aid\_\_median\_debt\_income\_0\_30000
16. aid\_\_median\_debt\_income\_30001\_75000
17. aid\_\_median\_debt\_income\_greater\_than\_75000
18. aid\_\_median\_debt\_independent\_students
19. aid\_\_median\_debt\_no\_pell\_grant
20. aid\_\_median\_debt\_noncompleters
21. aid\_\_median\_debt\_number\_completers
22. aid\_\_median\_debt\_number\_dependent\_students
23. aid\_\_median\_debt\_number\_income\_0\_30000
24. aid\_\_median\_debt\_number\_income\_30001\_75000
25. aid\_\_median\_debt\_number\_income\_greater\_than\_75000
26. aid\_\_median\_debt\_number\_independent\_students
27. aid\_\_median\_debt\_number\_overall
28. aid\_\_students\_with\_any\_loan
29. completion\_\_completion\_cohort\_less\_than\_4yr\_100nt
30. completion\_\_completion\_cohort\_less\_than\_4yr\_150nt
31. completion\_\_completion\_cohort\_less\_than\_4yr\_150nt\_pooled
32. completion\_\_completion\_rate\_4yr\_100nt
33. completion\_\_completion\_rate\_less\_than\_4yr\_100nt
34. completion\_\_completion\_rate\_less\_than\_4yr\_150nt
35. completion\_\_completion\_rate\_less\_than\_4yr\_150nt\_pooled
36. completion\_\_transfer\_rate\_4yr\_full\_time
37. completion\_\_transfer\_rate\_less\_than\_4yr\_full\_time
38. cost\_\_avg\_net\_price\_private
39. cost\_\_avg\_net\_price\_public
40. cost\_\_title\_iv\_private\_all
41. cost\_\_title\_iv\_public\_all
42. cost\_\_tuition\_in\_state
43. cost\_\_tuition\_out\_of\_state
44. school\_\_carnegie\_size\_setting
45. school\_\_carnegie\_undergrad
46. school\_\_degrees\_awarded\_highest
47. school\_\_degrees\_awarded\_predominant
48. school\_\_faculty\_salary
49. school\_\_ft\_faculty\_rate
50. school\_\_instructional\_expenditure\_per\_fte
51. school\_\_locale
52. school\_\_main\_campus
53. school\_\_men\_only
54. school\_\_minority\_serving\_aanipi
55. school\_\_minority\_serving\_annh
56. school\_\_minority\_serving\_hispanic
57. school\_\_minority\_serving\_historically\_black
58. school\_\_minority\_serving\_nant
59. school\_\_minority\_serving\_predominantly\_black
60. school\_\_minority\_serving\_tribal
61. school\_\_online\_only
62. school\_\_ownership
63. school\_\_religious\_affiliation

64. school\_\_state
65. school\_\_tuition\_revenue\_per\_fte
66. school\_\_women\_only
67. student\_\_demographics\_age\_entry
68. student\_\_demographics\_avg\_family\_income
69. student\_\_demographics\_avg\_family\_income\_independents
70. student\_\_demographics\_median\_family\_income
71. student\_\_family\_income\_dependent\_students
72. student\_\_family\_income\_independent\_students
73. student\_\_family\_income\_overall
74. student\_\_parents\_education\_level
75. student\_\_part\_time\_share
76. student\_\_retention\_rate\_four\_year\_full\_time
77. student\_\_retention\_rate\_four\_year\_part\_time
78. student\_\_retention\_rate\_lt\_four\_year\_full\_time
79. student\_\_retention\_rate\_lt\_four\_year\_part\_time
80. student\_\_size
81. student\_\_students\_with\_pell\_grant
82. student\_\_valid\_dependency\_status
83. repayment\_rate

### 3. Linear Regression Model

Since our data mining problem is a predictive one, a prediction model was created for the data. The data to be predicted i.e. repayment\_rate is a continuous attribute hence, regression was found to be the best technique to handle continuous data. Linear regression was chosen as the technique. The output of the model is the repayment\_rate for the test values given.

Linear regression is a linear model, e.g. a model that assumes a linear relationship between the input variables (x) and the single output variable (y). More specifically, that y can be calculated from a linear combination of the input variables (x).

Learning a linear regression model means estimating the values of the coefficients used in the representation with the data that we have available.

Given the representation is a linear equation, making predictions is as simple as solving the equation for a specific set of inputs.

Let's make this concrete with an example. Imagine we are predicting weight (y) from height (x). Our linear regression model representation for this problem would be:

$$y = B_0 + B_1 * x_1$$

or

$$\text{weight} = B_0 + B_1 * \text{height}$$

Where  $B_0$  is the bias coefficient and  $B_1$  is the coefficient for the height column. We use a learning technique to find a good set of coefficient values. Once found, we can plug in different height values to predict the weight.

In case of the model, the repayment\_rate is given as the linear combination of the chosen 83 attributes as:

`repayment_rate = -0.7232 * admissions__act_scores_midpoint_english + 1.1902 * admissions__act_scores_midpoint_math + 0.5365 * admissions__admission_rate_overall + 0.719 * admissions__sat_scores_average_overall -0.4737 * admissions__sat_scores_midpoint_math + ...and so on.`

edx Username: roseenah1

In context of Azure ML Studio, the model is a simple linear regression. There are a few types of linear regression modules, the one named "Linear Regression", found under "Regression" was used.

The "Learning Rate" parameter allows us to set how much difference we see from tree to tree. This describes quite well as "the learning rate determines how fast or slow the learner converges on the optimal solution. If the step size is too big, it might overshoot the optimal solution. If the step size is too small, training takes longer to converge on the best solution." The "Random Number Seed" parameter allows us to create reproducible results for presentation/demonstration purposes. The model was tested before selecting it as final Model. For testing, Model was trained with 60% of the combination of train\_values and train\_labels, and tested with the remaining 40% yields following results after the split:

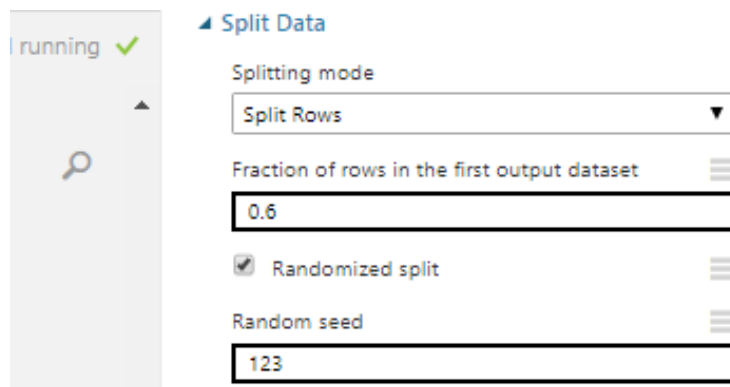


Fig. Splitting of data

The data after Score Label step is:

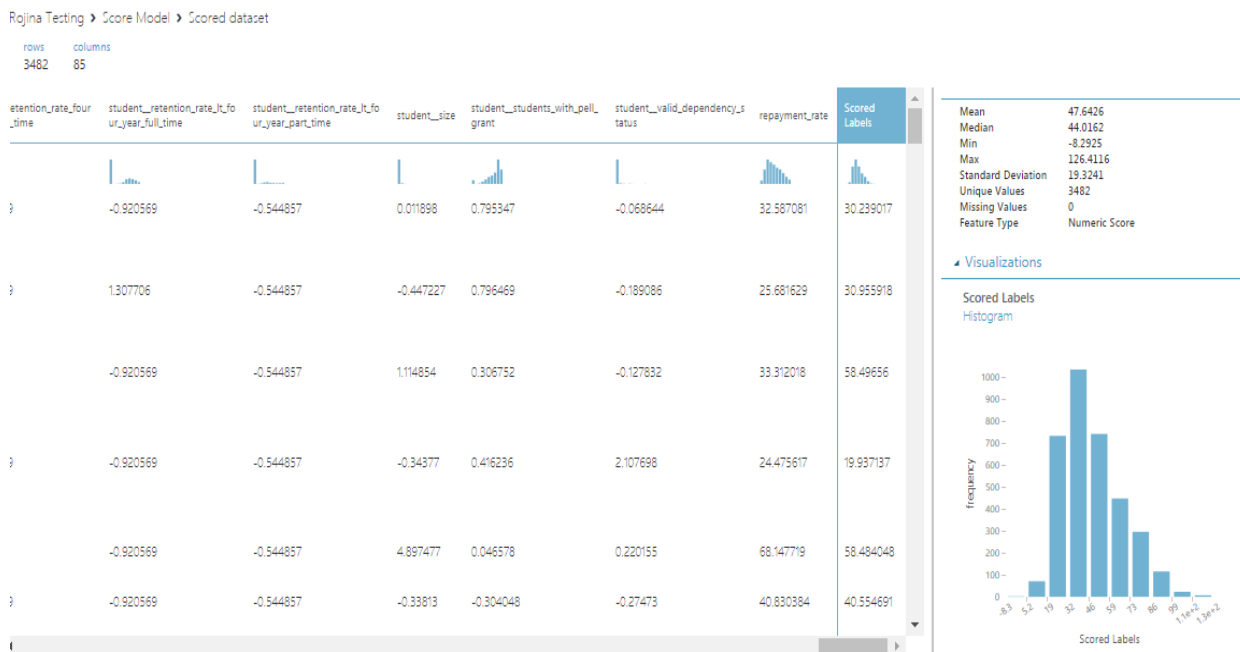


Fig. Visualization of scored dataset

edx Username: roseenah1

Mean Absolute Error: 5.528142

Root Mean Squared Error: 7.421821

Relative Absolute Error: 0.352459

Relative Squared Error: 0.120887

Coefficient of Determination: 0.903207

The module with the highest value for Coefficient of Determination is taken as best. And, here the Coefficient of Determination is near to 0.1 so it is taken as significant. Hence, the model was chosen as final regression model.

#### Error Histogram

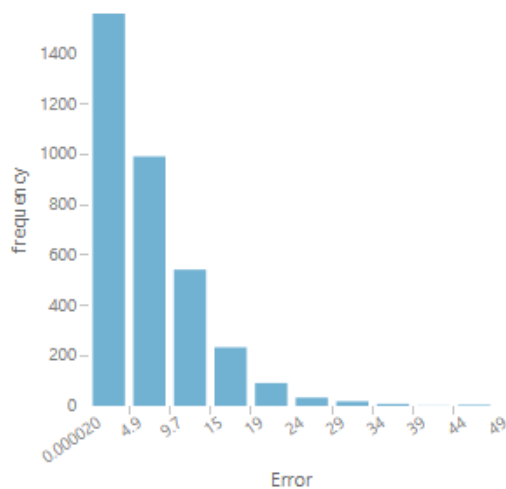


Fig. Error histogram of the linear regression model

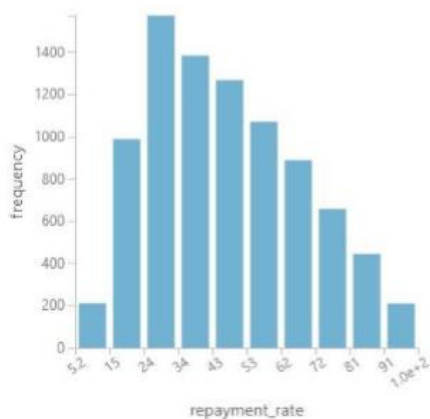


Fig.Repayment\_rate acutal data histogram

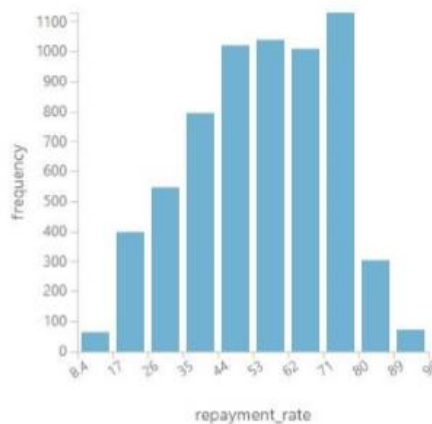


Fig. Repayment\_rate predicted histogram

#### 4. Conclusion

Thus, going through the critical data mining techniques, the required value with Root Mean Squared Error (RMSE) of 9.492911 was predicted. The linear regression model was best suited for the data, upon individual evaluation and the error obtained is acceptable. Hence, the model was concluded with the above procedures.