

# College Motivation: Predicting Students Going to College

## By: Rebekah Sander

### ABSTRACT

The goal of this project was to use predictive analytics to identify whether a student will make the decision to go to college. This "college motivation" dataset collected 1000 high school students in Mexico who graduated in 2014. A regression was used to make predictions of a student going to college based on variables that have to do with their academic and personal backgrounds.

### METHODS—Preparing the Data

- 1.) Checking for missing data:
  - There is no missing data
- 2.) Recoding Variables:
  - "will\_go\_to\_college" and "parent\_was\_in\_college": Originally character values "True" and "False", changed to 1 and 0 for binary.
- 3.) Data Split:
  - 6:2:2 data split ratio.
  - 60% of data into new dataframe, "college\_train", 20% into new dataframe "college\_validation", 20% into new dataframe, "college\_test".

### METHODS—Building the Best Model

Regression Model:

- Logistic regression model based on data frame, "college\_train" with will\_go\_to\_college as the binary dependent variable, all other variables are used as independent variables.

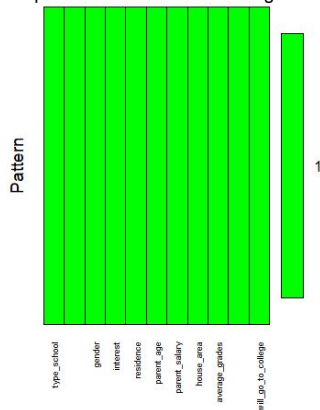
Choosing the best threshold for binary prediction:

- Best threshold determined through "college\_validation" dataframe.

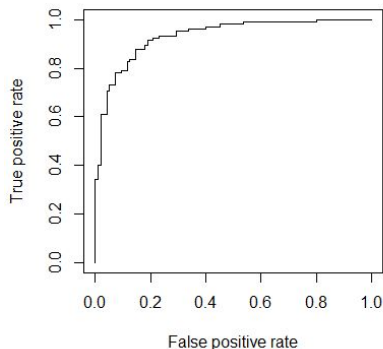
Performance and Evaluation Metrics

- PMSE
- AUC-ROC
- Accuracy

Graph 1: Visualization of Missing Data



Graph 2: AUC-ROC curve



### DESCRIPTION

**will\_go\_to\_college:** Student goes to college  
**type\_school:** Type of high school the student went to  
**school\_accreditation:** Student's high school's accreditation  
**gender:** Gender of the student  
**interest:** Student's interest in going to college  
**residence:** Where the student lives  
**parent\_age:** Age of student's parent  
**parent\_salary:** How much the student's parent makes  
**house\_area:** Area of student's house  
**average\_grades:** Average grades of student  
**parent\_was\_in\_college:** Student's parent went to college

### RESULTS

		REFERENCE	
		0	1
PREDICTION	0	86	22
	1	9	83

Best Threshold = 0.64  
Accuracy = 0.86  
PMSE = 0.33  
AUC-ROC = 0.92

### CONCLUSIONS

During this predictive analysis, a logistic regression model was utilized to identify high school students in Mexico who are likely to attend college based on their personal and academic life. The primary objective was to pinpoint and be able to reach out to students who may not have plans to pursue higher education.

The model achieved an AUC-ROC evaluation metric of 0.92. This high value indicates the model's strong ability to correctly predict whether a student is college-bound. However, when it comes to generalizability, it is important to acknowledge the dataset is based in Mexico. This model may not be the case when looking at different regions. Since this dataset was a collection of students in Mexico, this model may not be extended to other countries. To improve generalizability, more data could be collected from various other regions and countries.

Overall, this analysis provides valuable insights into identifying students at risk of not attending college so that intervention may occur.