# INDIVIDUAL ASSIGNMENT
## RESPONSE TEMPLATE

Please fill in your name and submission date below:

〉 Ckalib S. Nelson

〉 May 15th, 2021

Please fill-in your answers to your questions below. Feel free to incorporate any graphs or other information you may consider. I have included a recommended length of answer for orientation. Remember your answers to these questions will weight 70% of your final grade on the Individual Assignment.
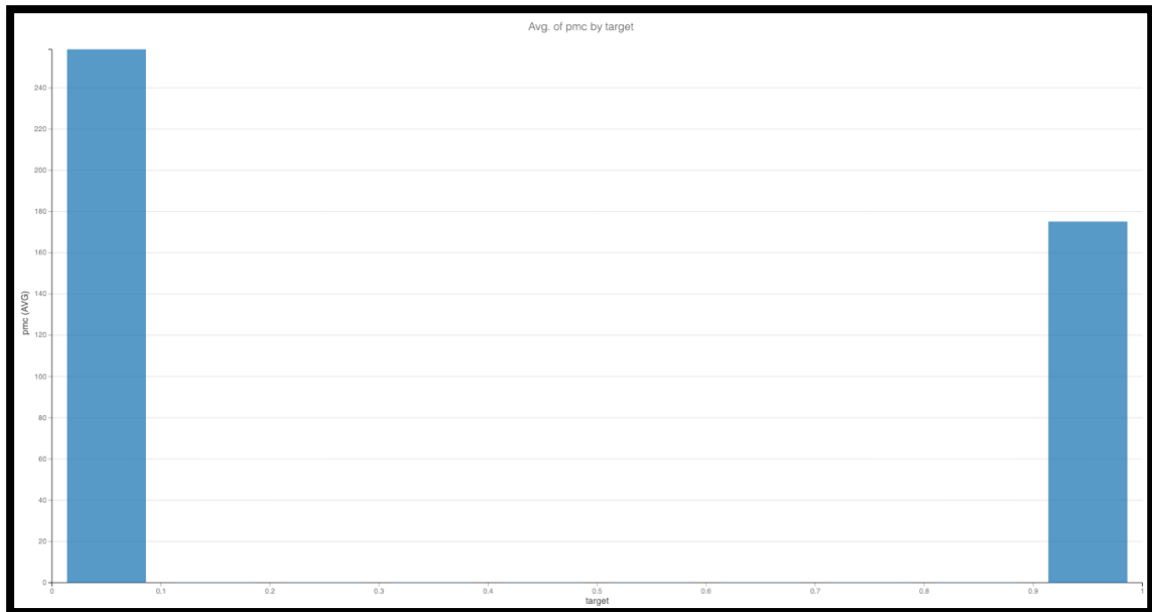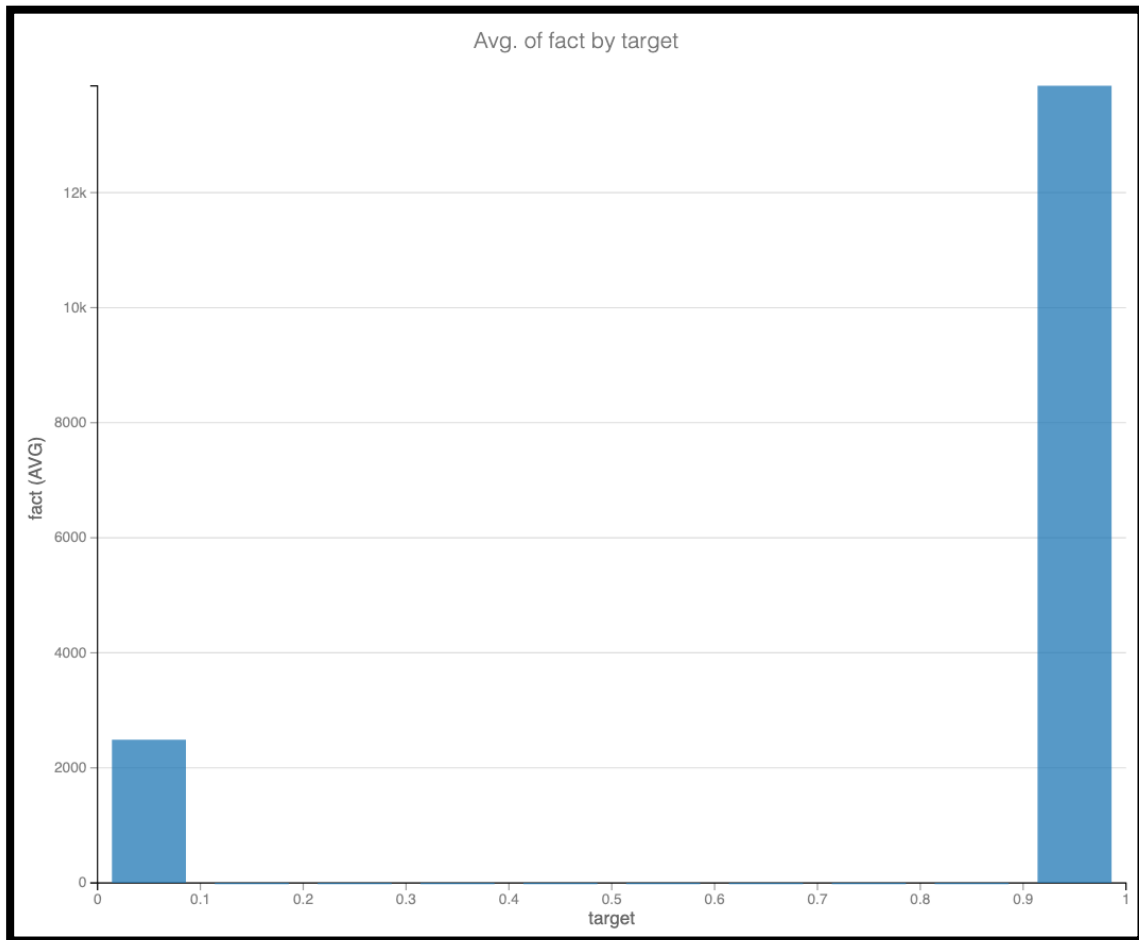
## Table of Contents

# QUESTION 1.

## Understanding the Data [recommended max. 150 words]

How does the target variable behave when compared to the other variables? What sort of interactions do you see between the variables? What sort of correlation? Which are the most correlated variables? Do you see a business-sense?

Here are a few visualizations that highlight the relationship between the target variable and the independent variables:



Avg. of pmc by target

Per the chart above, illustrating the average collection period (**pmc**) by **target**, the higher the **pmc** the less likely the company will contract the POS in the following 3 months. Thus, there is a negative correlation between **pmc** and the **target**. In a business sense, this makes sense.

Avg. of fact by target

Per the chart above, illustrating the average amount of money the cards emitted by the bank spent in the POS owned by the company (**fact**) by **target**, the higher the fact the more likely the company will contract the POS in the following 3 months. Thus, there is a positive correlation between **fact** and the **target**. In a business sense, this makes sense.

Avg. of noprotra by target

Per the chart above, illustrating the number of services the company has contracted with the bank (**noprotra**) by target, the higher the **noprotra** the more likely the company will contract the POS in the following 3 months. Thus, there is a positive correlation between **noprotra** and the **target**. In a business sense, this makes sense.
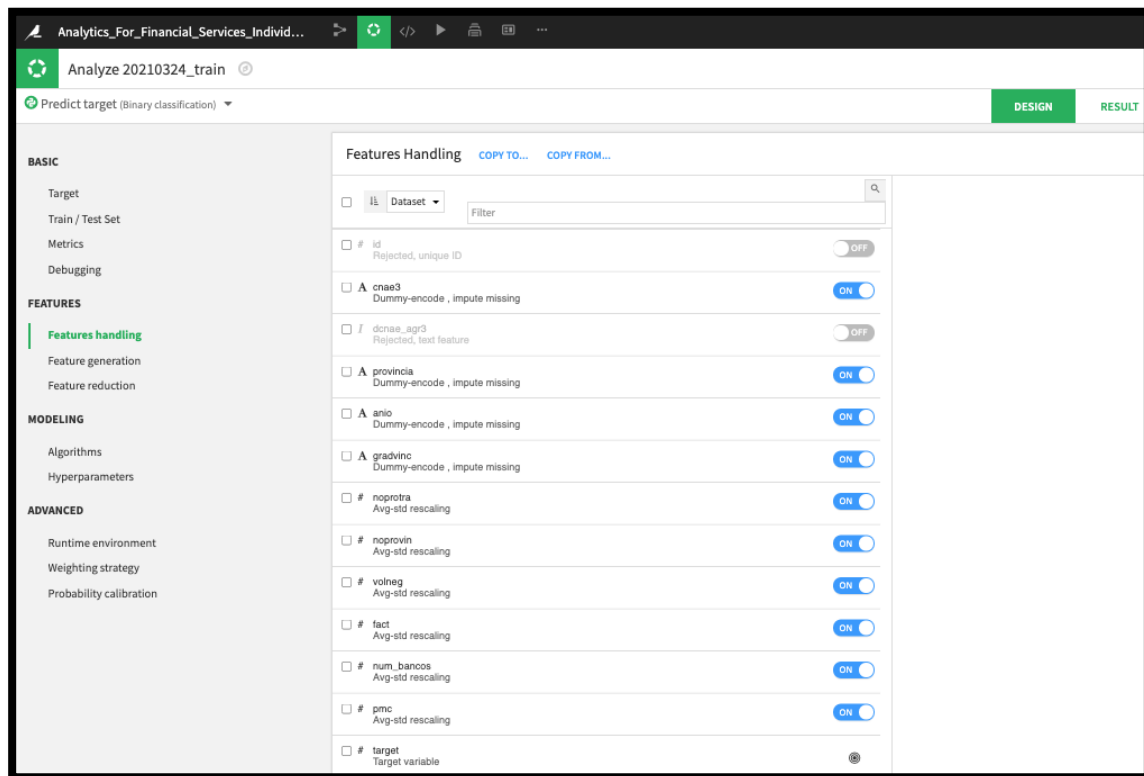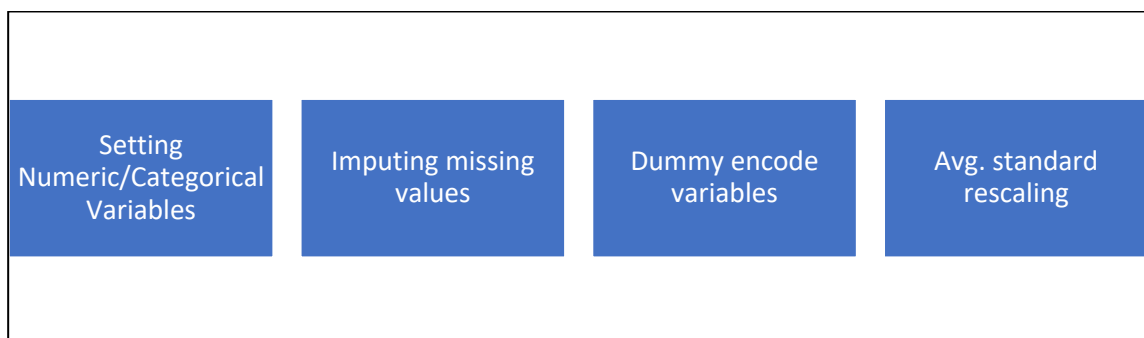
## QUESTION 2.

### Preparing the Data [recommended max. 100 words]
Would you recommend doing any sort of cleaning or transforming the variables before working on the model? (eg. missing values, outliers, etc…)

Yes, I would recommend doing a bit of cleaning and transforming of the variables before working on the model given that *a few columns have missing values and are not of the necessary data type*:



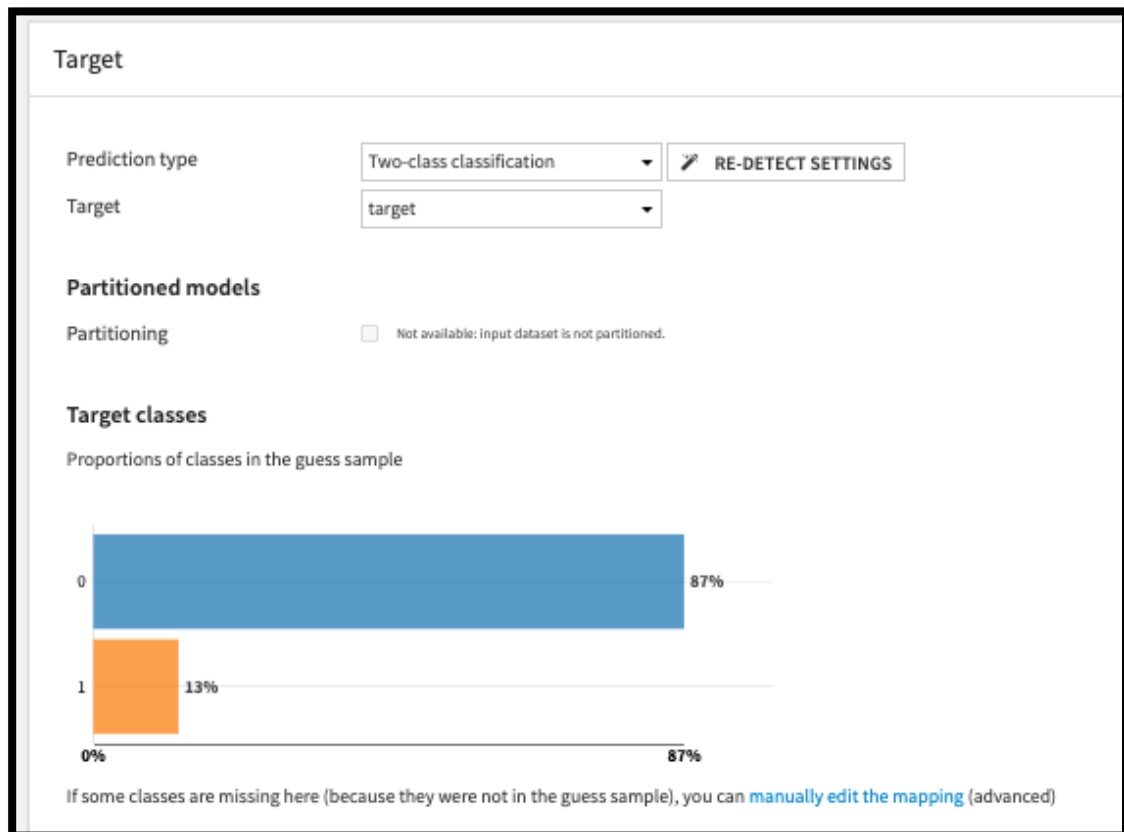The cleaning and transforming needed can be summarized as follows:

| Setting Numeric/Categorical Variables | Imputing missing values | Dummy encode variables | Avg. standard rescaling |
|---|---|---|---|

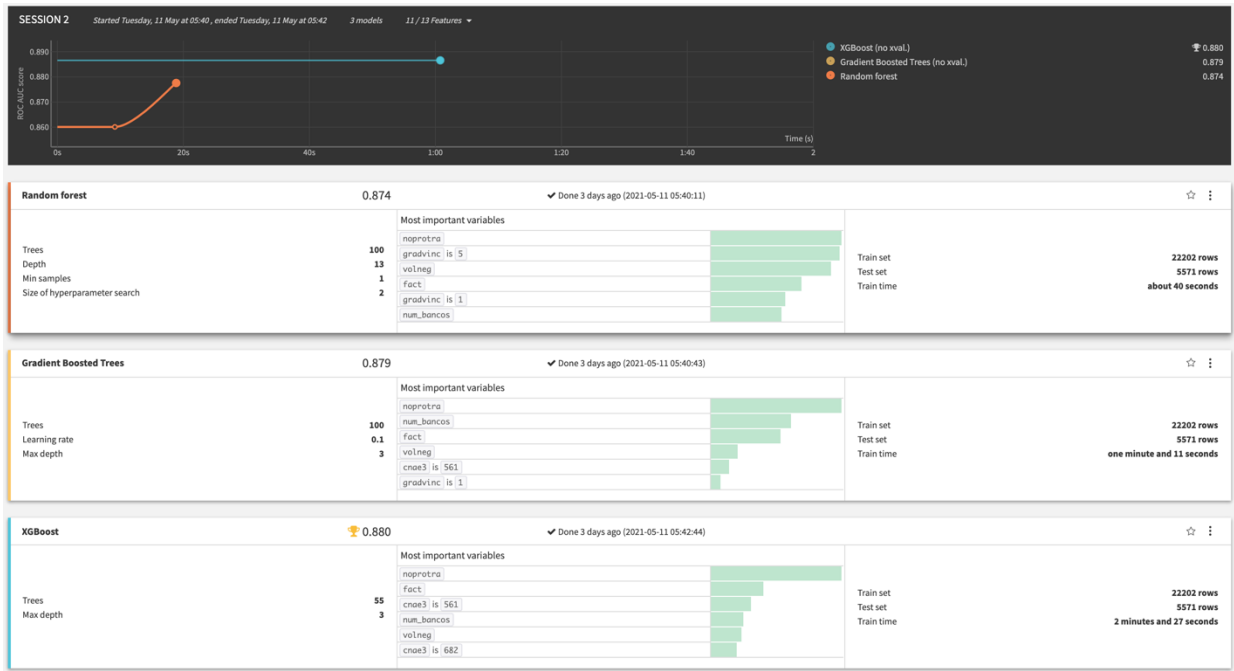# QUESTION 3.

## The Model [recommended max. 150 words]
What type of model have you used (ie. Supervised/unsupervised, binary...) and why?

Given that we do know the target variable, we will have to rely upon a supervised model. Furthermore, the target variable is binary; thus, we will have to leverage a supervised classification model.

After a bit of trial and error, I settled on the XGBoost algorithm, or the e**X**treme **G**radient **B**oosting algorithm. This particular algorithm performed better (per AUC) than its competitors:
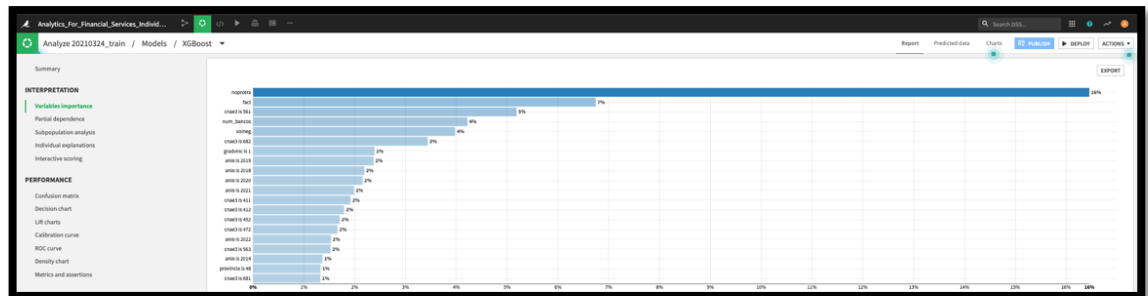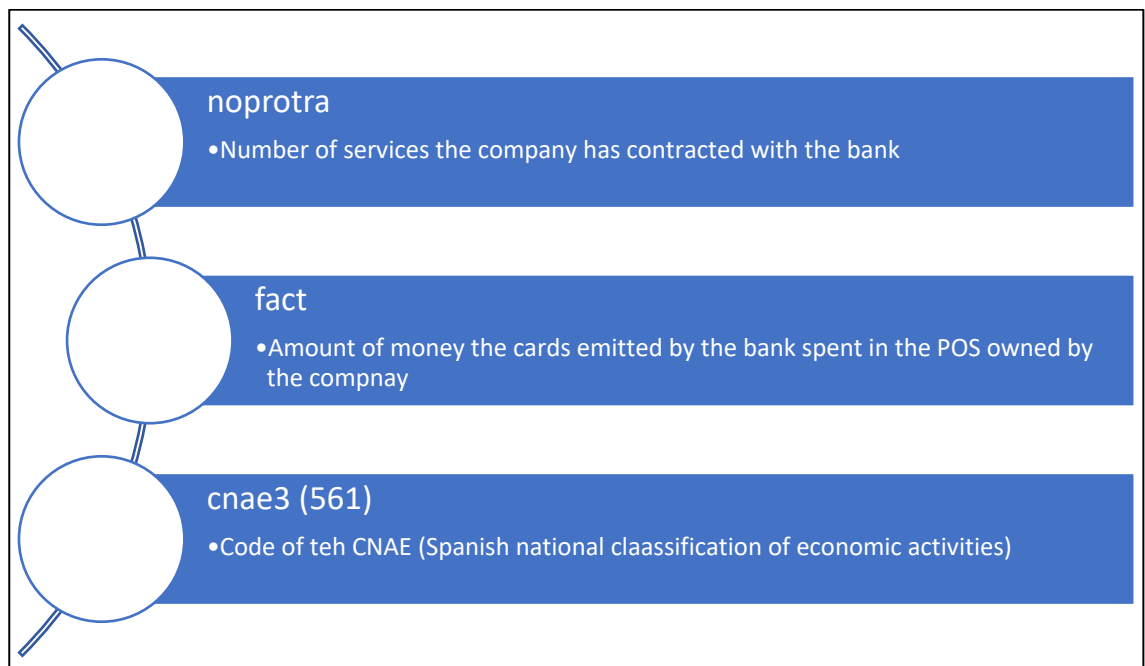
# QUESTION 4.

## The Outcome [recommended max. 150 words]
What are the most significative variables in the model?

The most significant variables in the model include the following:



Here is a list of the top three variables of importance:



**noprotra**
- Number of services the company has contracted with the bank

**fact**
- Amount of money the cards emitted by the bank spent in the POS owned by the compnay

**cnae3 (561)**
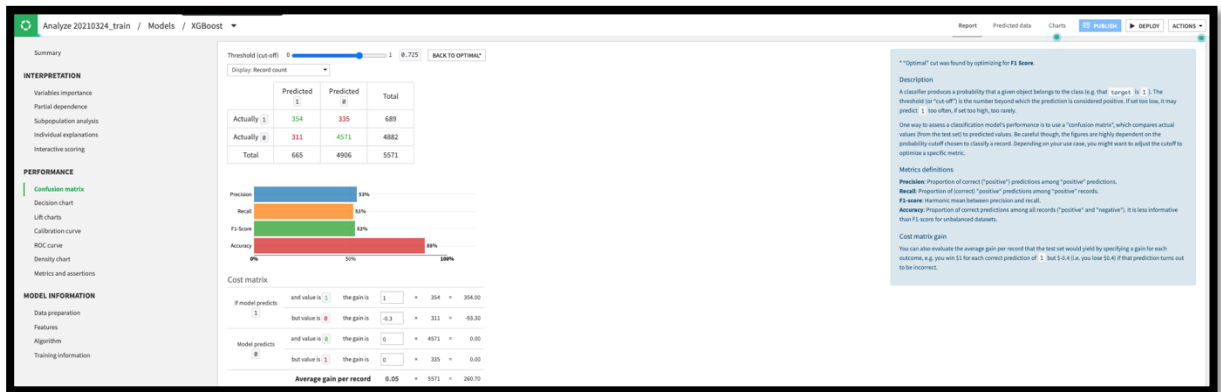- Code of teh CNAE (Spanish national claassification of economic activities)

Unsurprisingly, the correlations we uncovered during the exploratory data analysis rings true when evaluating the model. For example, the higher the **noprotra** the more likely the company will contract the POS in the following 3 months. In other words, companies are much more likely to purchase when the number of services the company has contracted with the bank is high. Similarly, there is a positive correlation between the target and **fact**. Essentially, companies are much more likely to purchase when the amount of money the cards emitted by the bank spent in the POS owned by the company is high.

**Please complete the Data-set-test variables with the outcome of your model –** remember the accuracy of the model will weight 30% of your final grade in the individual assignment.  Please Indicate here the accuracy of the model:

Using the threshold (cut-off) of .725, the model has the following confusion matrix and performance metrics:



In regard to the business case, the worse outcome from this model is informing the audience that the model predicts a company would not buy although in reality the company does buy. In other words, we want to reduce the number of **False Negatives,** or the aggregate occurrences the model predictes "No, the company won't buy" although it was actually "Yes, the company will buy." Thus, if we were to reduce the number of False Negatives, we would want the **threshold** to be a bit more generous, or lower than it currently is. However, if we want to optimize for accuracy, we will want to keep the threshold where it is.

*Example Answer: The accuracy of the model in the test data **is 0.75**. The model predicts well **2.139 cases** over the total 2.852 cases of the test dataset (**77%** of the cases - accuracy).*

**Your Answer:** The accuracy of the model in the test data **is .884.** The model predicts well **2521 cases** over the total 2.852 cases of the test dataset (**88%** of the cases - accuracy).
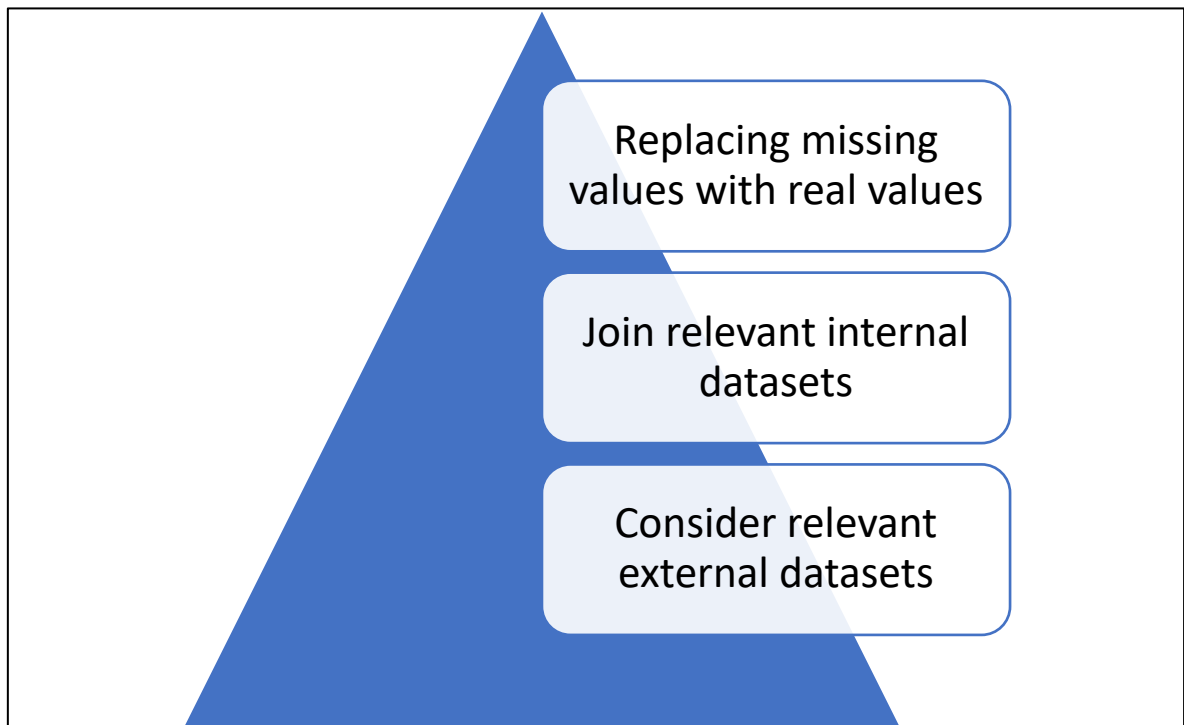
## QUESTION 5.

**Improving the model** [recommended max. 100 words]

What type of improvements to the Data-set or the steps followed through the analysis would you implement to <u>improve the effectivity of the model?</u>
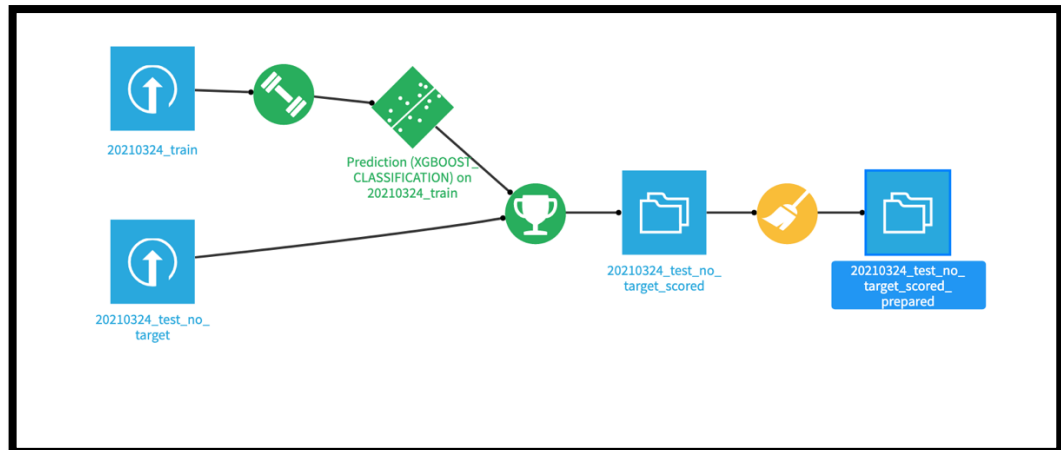
I would suggest the following improvements to improve the effectivity of the model:



Replacing missing values with real values

Join relevant internal datasets

Consider relevant external datasets

## ANY OTHER COMMENTS [recommended max. 100 words]

Here is a snippet of the "flow," or the process by which this analysis was conducted, per Dataiku:



In summary, I completed the following machine learning process to build predictions for the test dataset: