

phData Case Study

(Business Version)

Michael Rowlands

Tax Company's Challenge

For years, Tax Company has been unable to identify leads that will result in sales of your software.

This has resulted in thousands of dollars in unnecessary labor and advertising costs.

Tax Company has created a dataset containing two years of customer information and if you were successful at selling to each customer.

phData's Solution

Investigate the dataset, and determine if a machine learning approach is viable.

If so, create a model to predict if a lead will convert.

Preliminary Assumptions

The dataset is historical so we assume future data is similar to the historical data.

Customers only appear once in the dataset.

Tax Company has no prior knowledge on what features are predictive of sale.

There is a fixed profit for sale (true positives) and fixed loss for attempted sales (false positives).

Exploratory Analysis

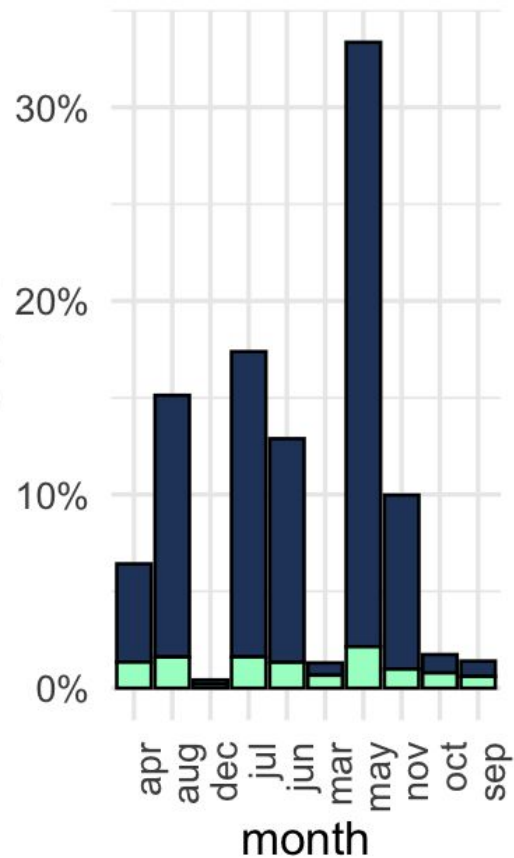
Basic Dataset Info

Full dataset contains 41,188 sale attempts and 22 variables of customer information along with if the sale attempt was successful

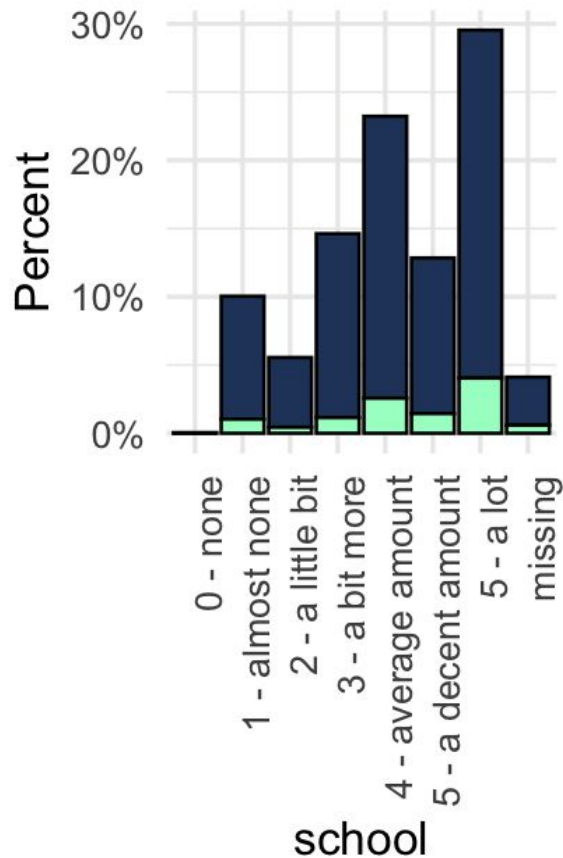
Dataset contains information like customer age, employment type, marriage status, and a lot of variables with undescriptive names like “c3”.

We put aside 20% of the dataset to use as a benchmark for how the model we build will perform on new data

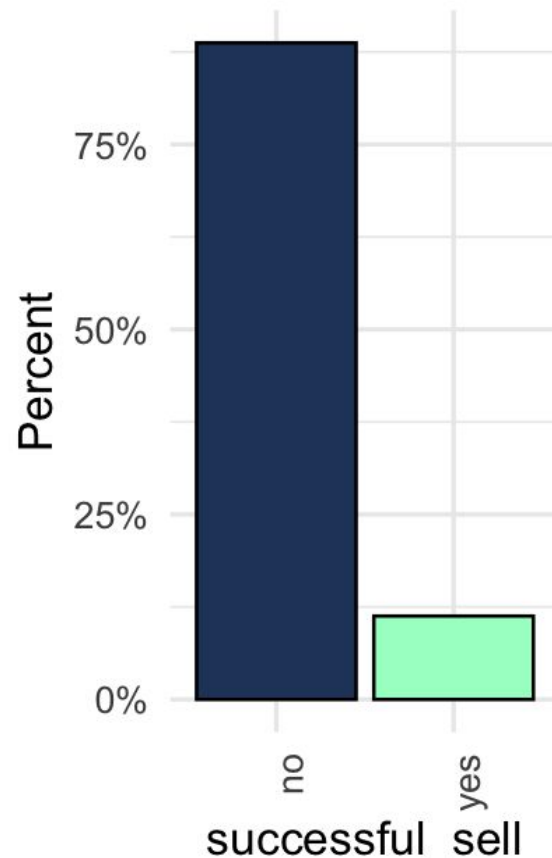
month



school

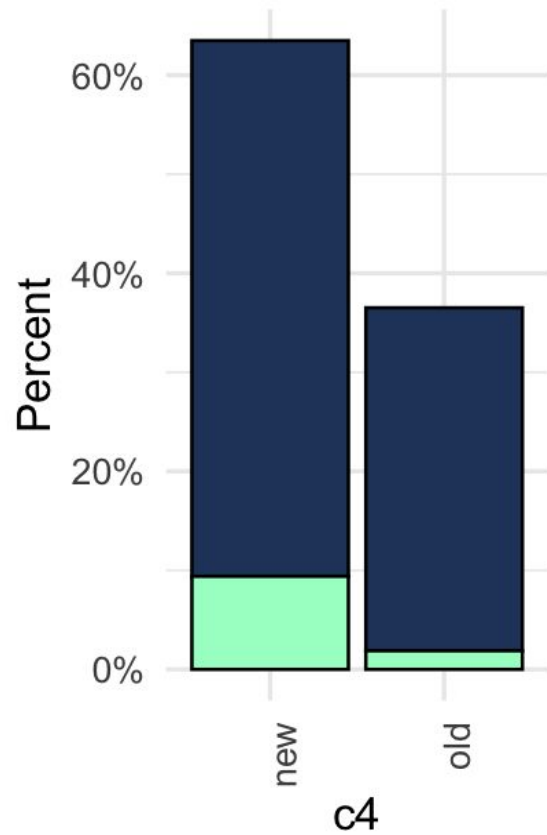


successful_sell

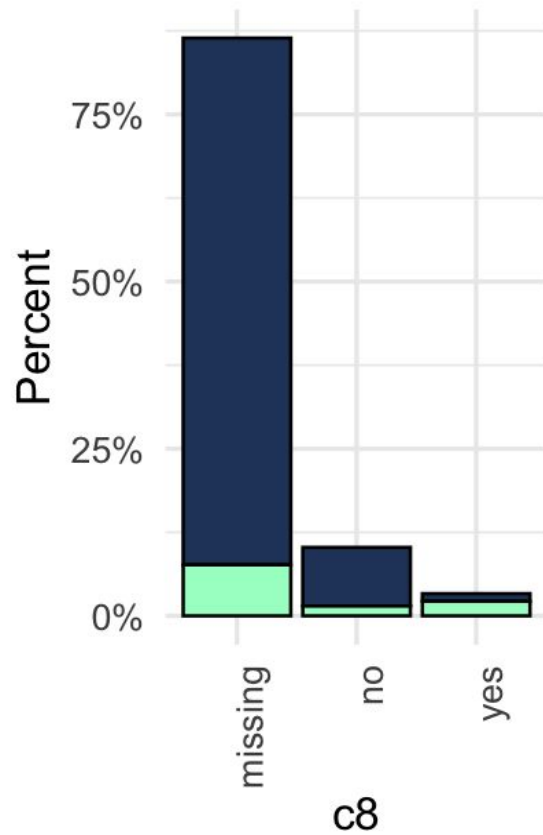


successful_sell no yes

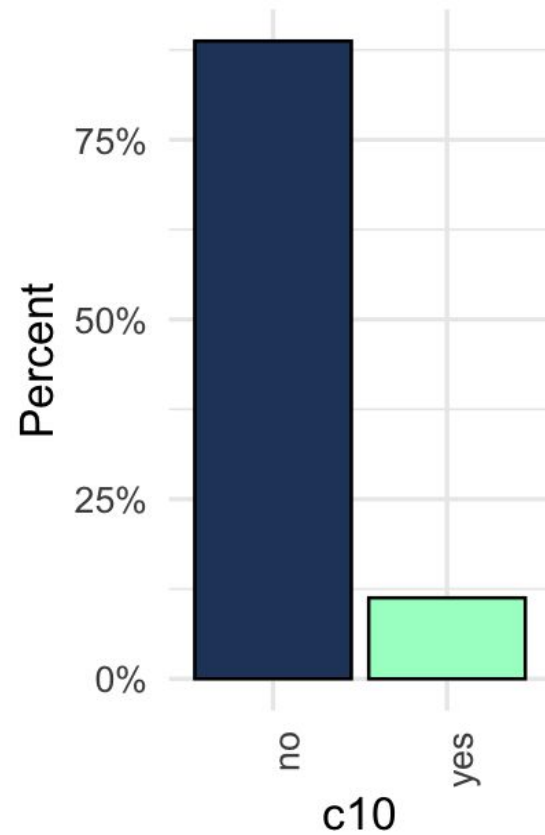
c4



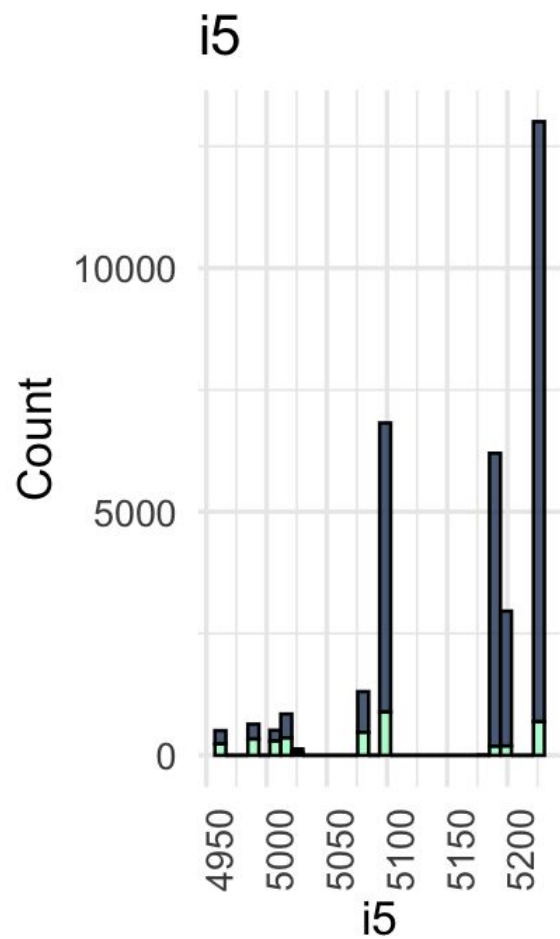
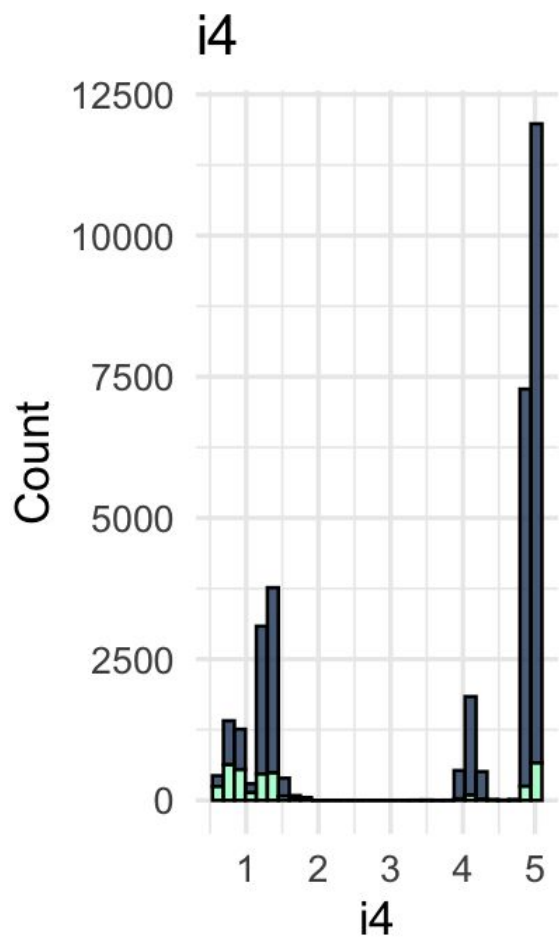
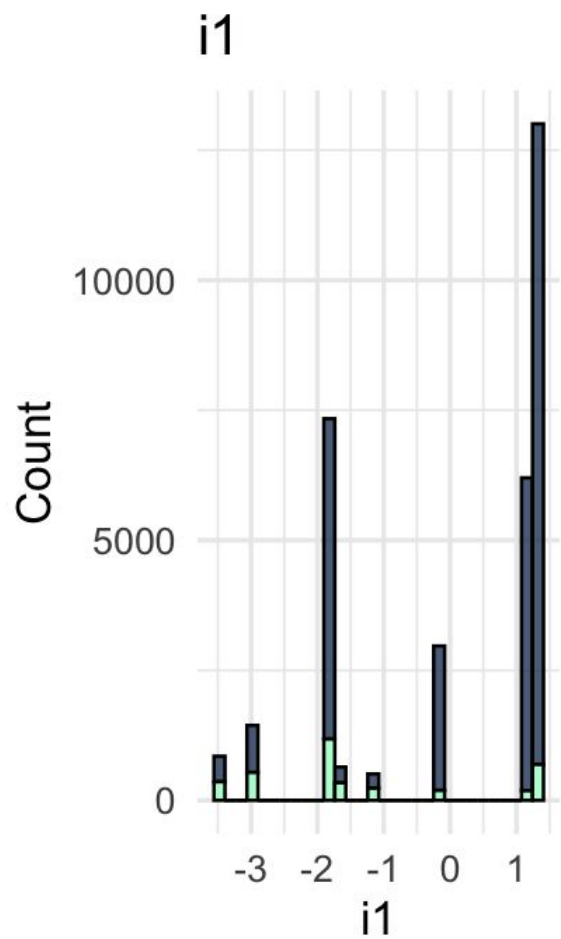
c8



c10

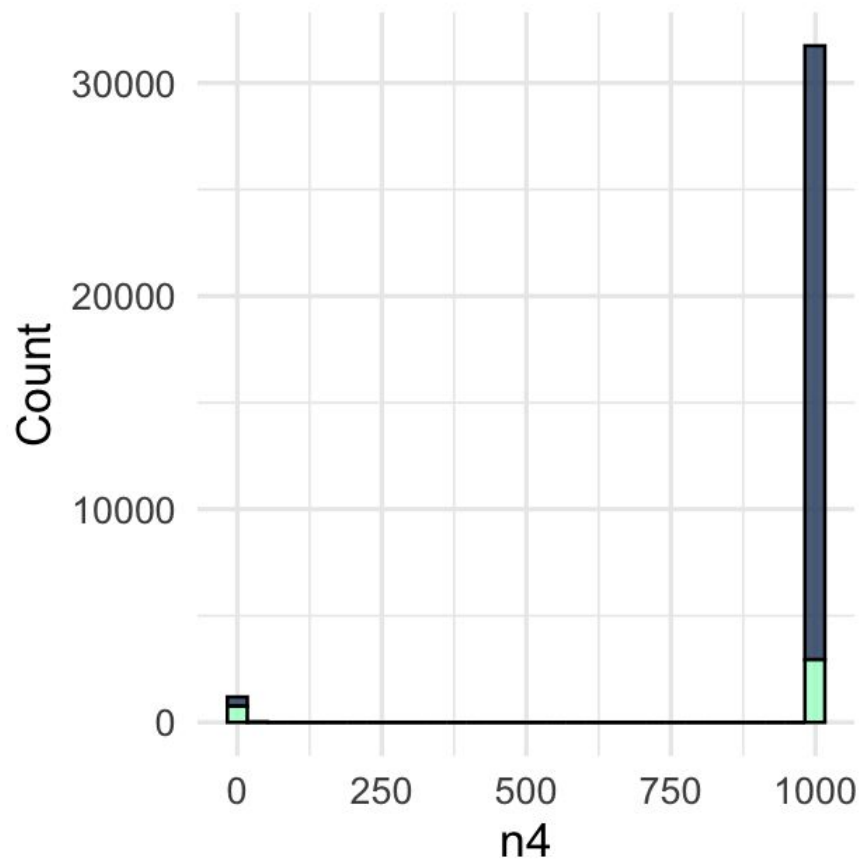


successful_sell no yes

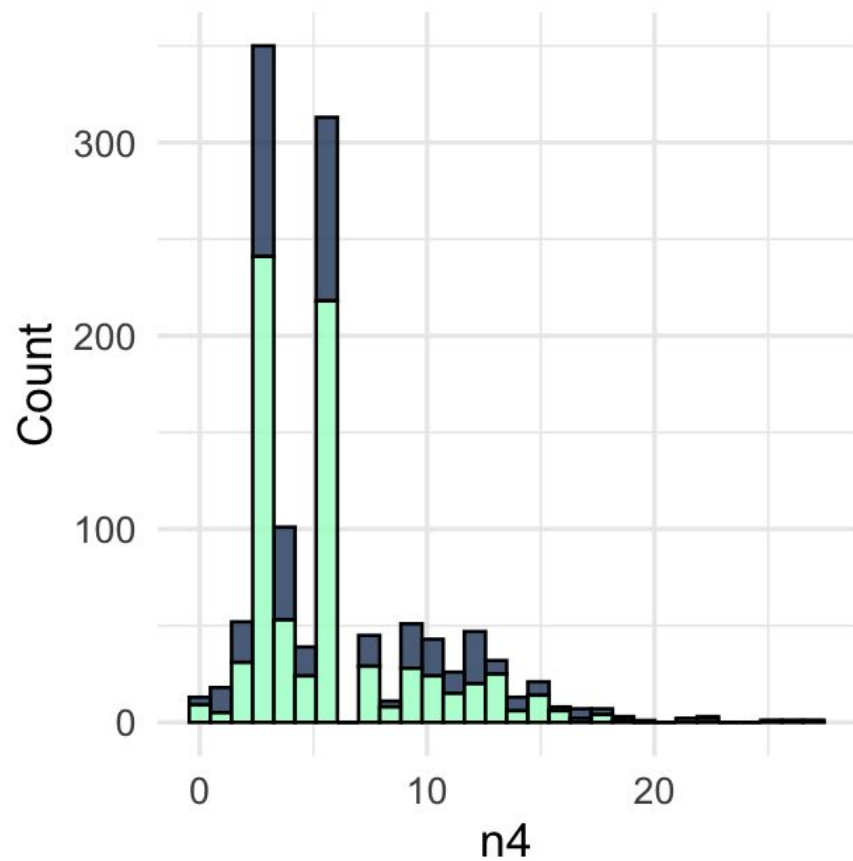


successful_sell no yes

n4

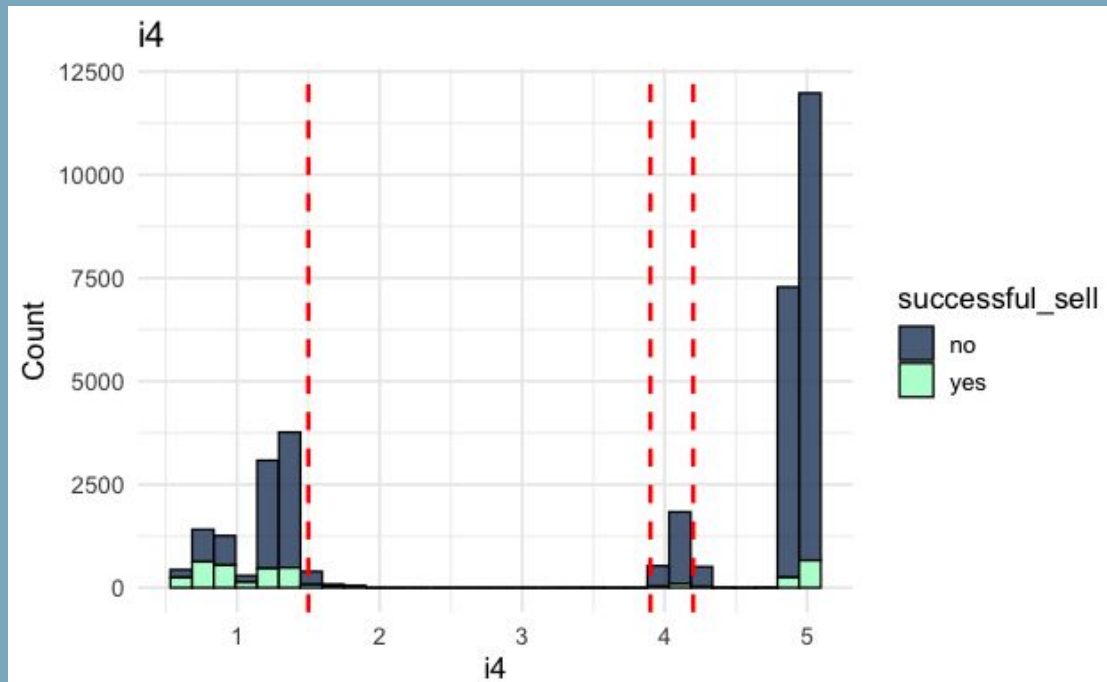


n4 Filtered



successful_sell no yes

We prepared the data to work well in a variety of different types of models



Modeling and Results

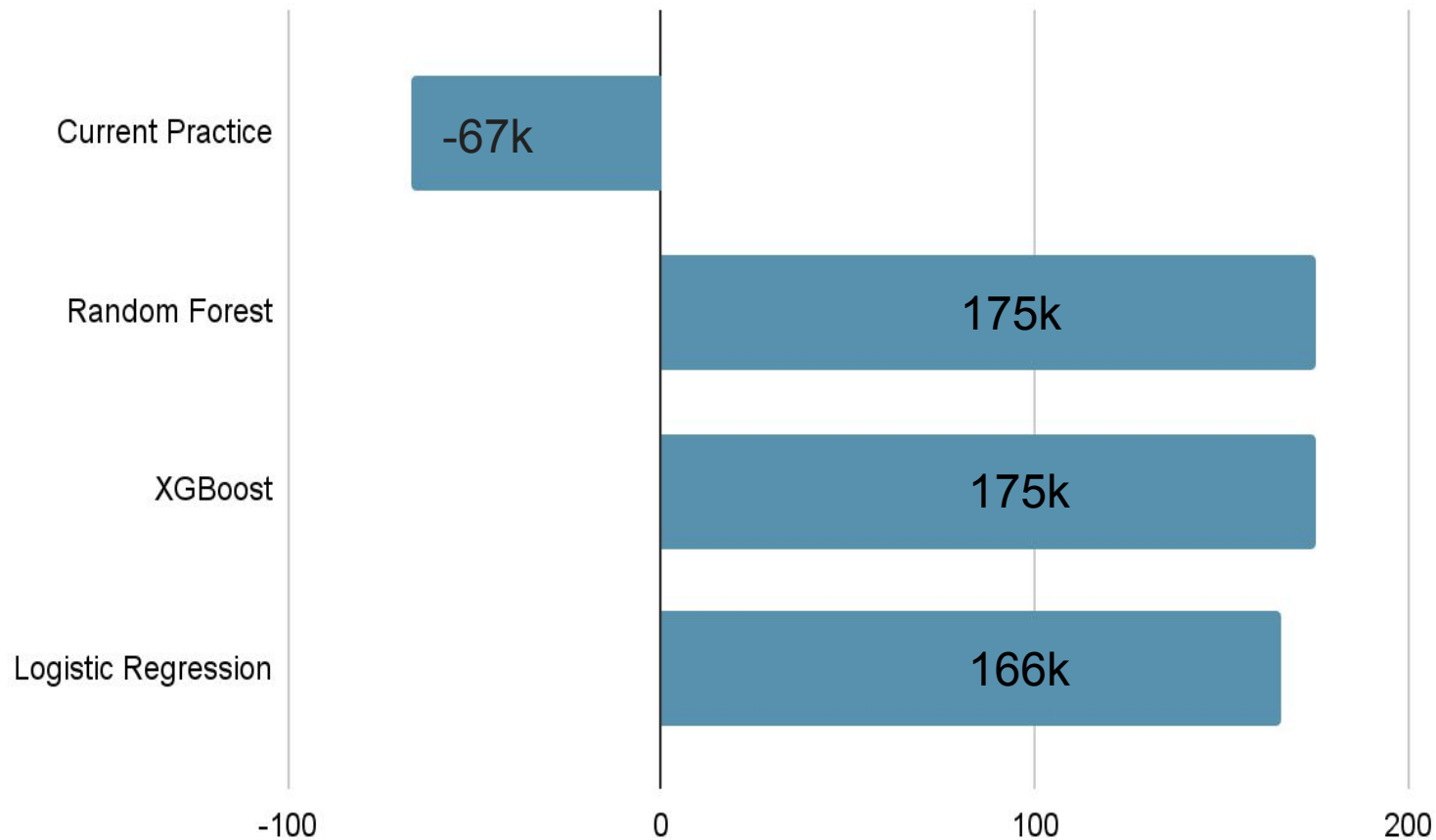
We evaluated the performance of three different models along with a “naive” model which captures the current operating practice of Tax Company

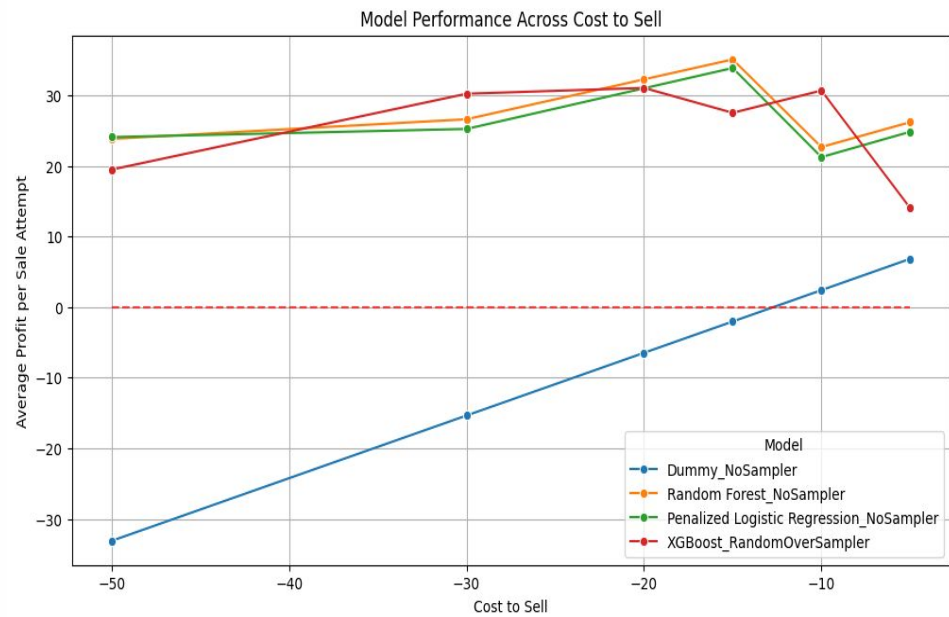
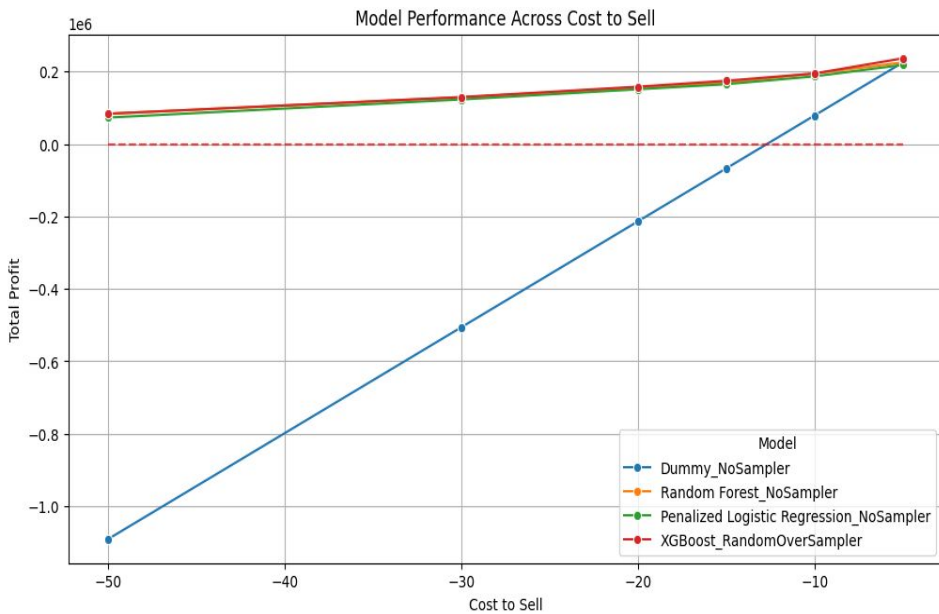
Performance of each model was optimized to maximize profit according to your current profit per sale and the cost of an unsuccessful attempt

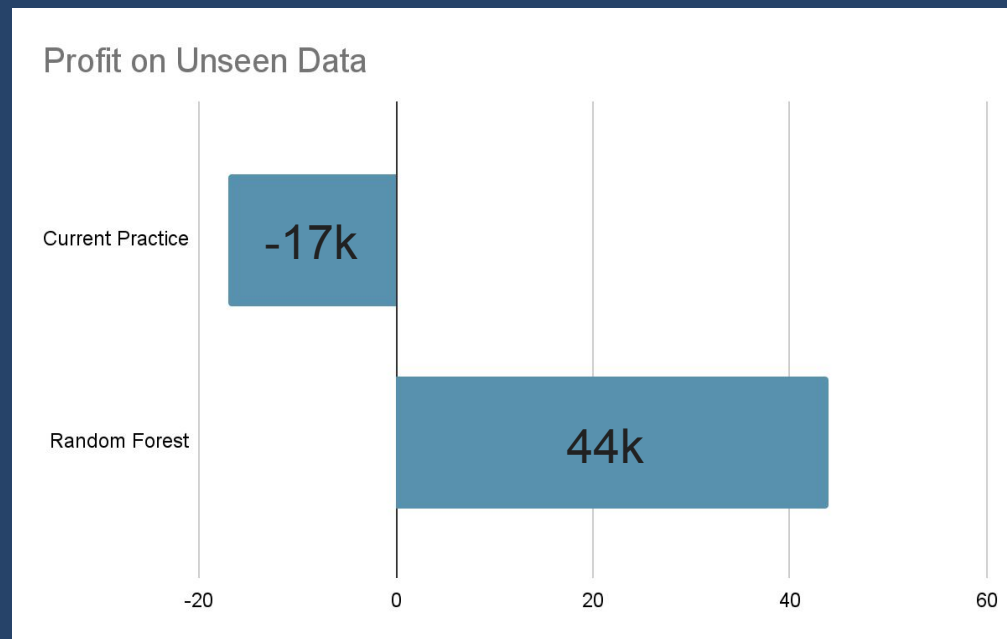
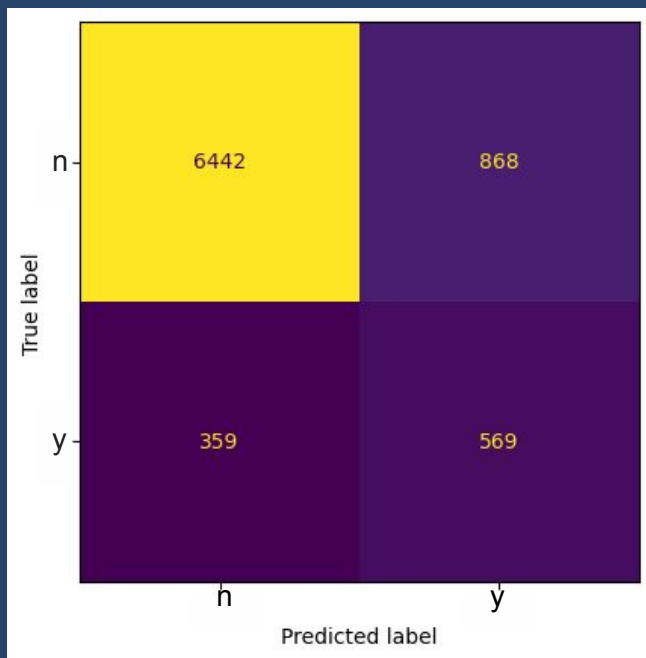
The models we used were:

- 1) Logistic Regression
- 2) Random Forest
- 3) XGBoost

Profit on Training Set





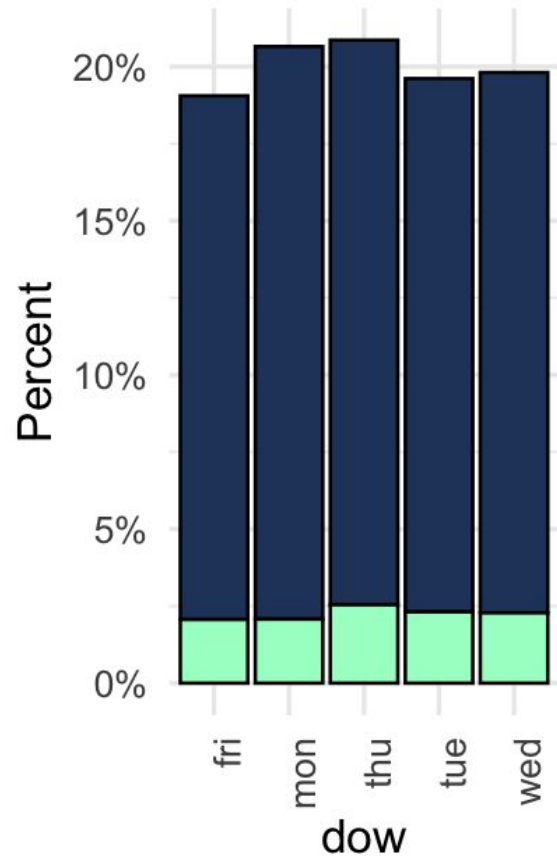


Next Steps and Places for Improvement

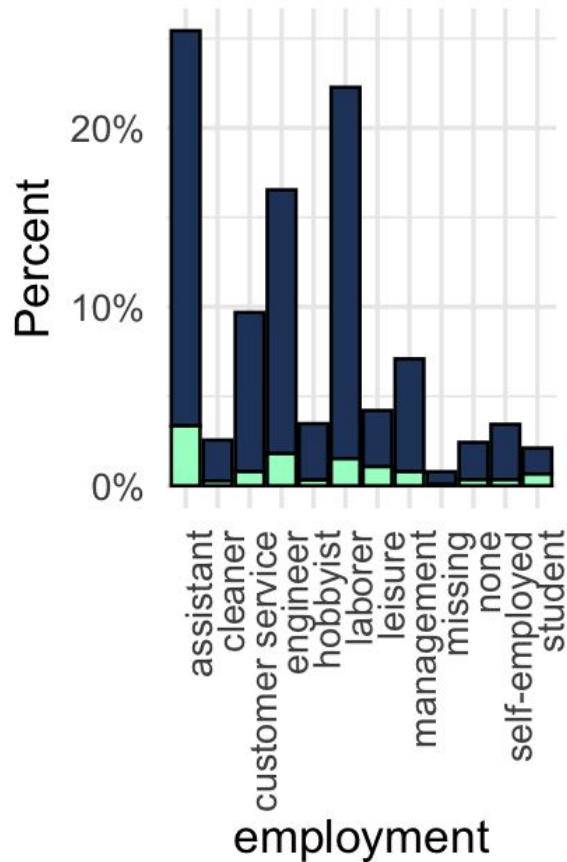
- Deploy and utilize the model so you can use it to generate likely leads
 - phData will monitor performance to ensure likely leads are identified
- Look to utilize domain knowledge and other data sources if available in future modeling
 - If we can leverage strong domain knowledge, we can construct our models to be even better

Questions???

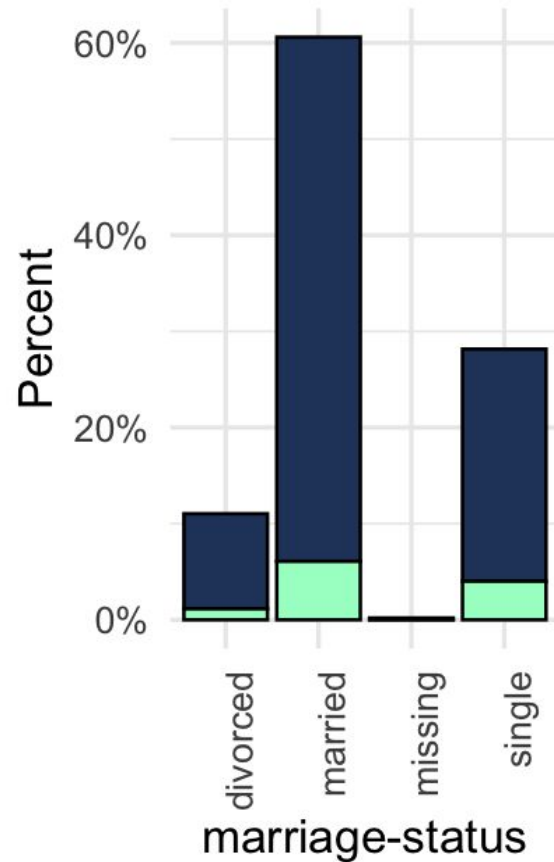
dow



employment

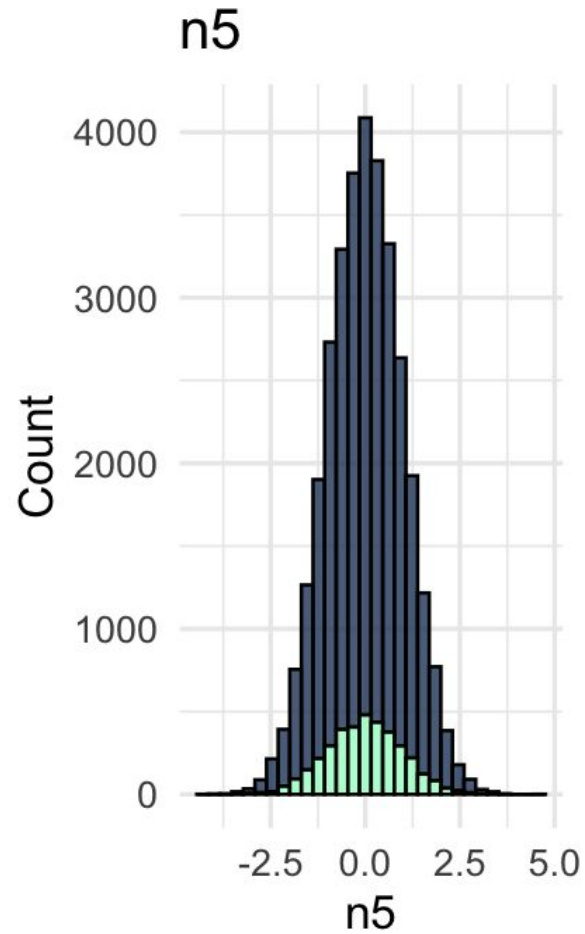
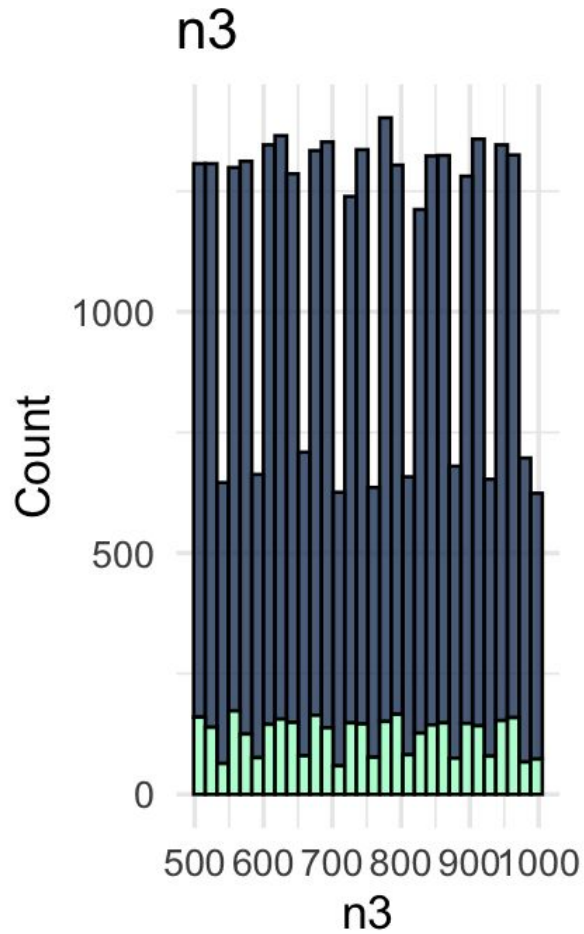
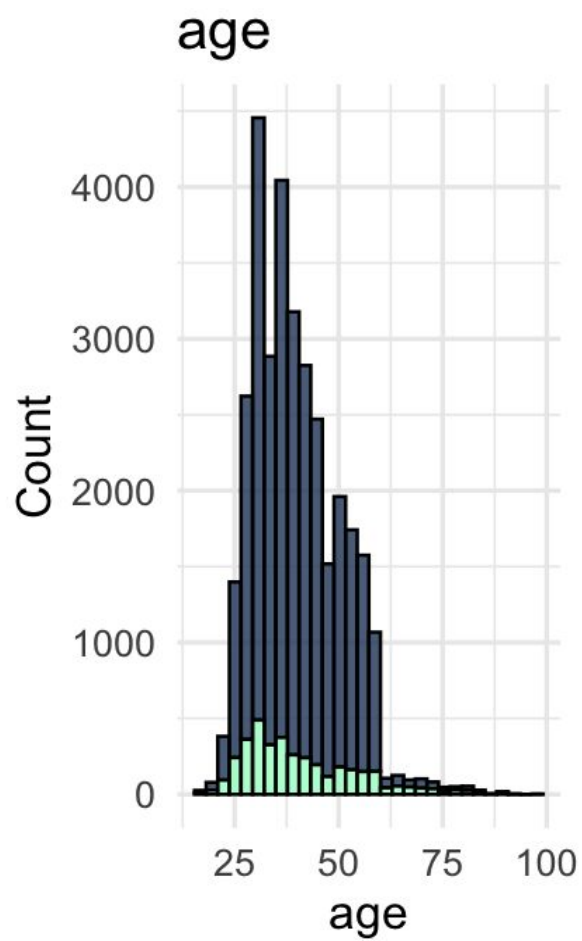


marriage-status



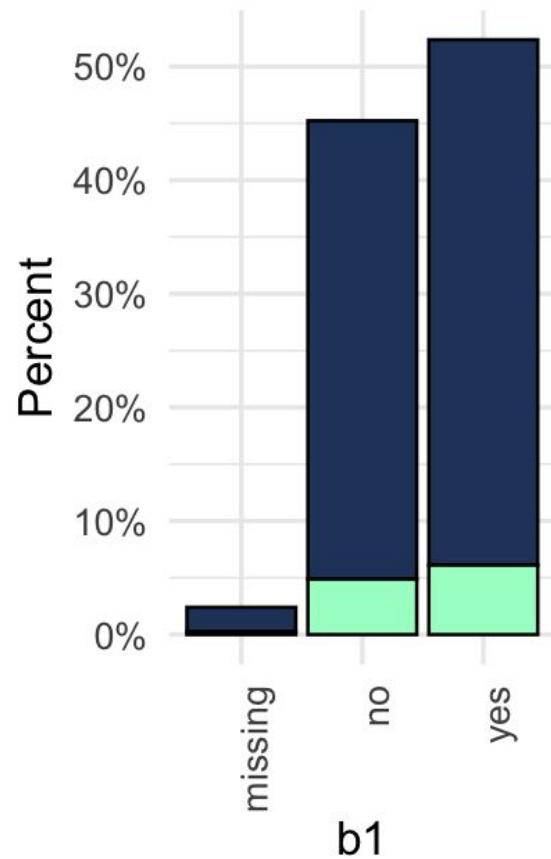
successful_sell

no	yes
----	-----

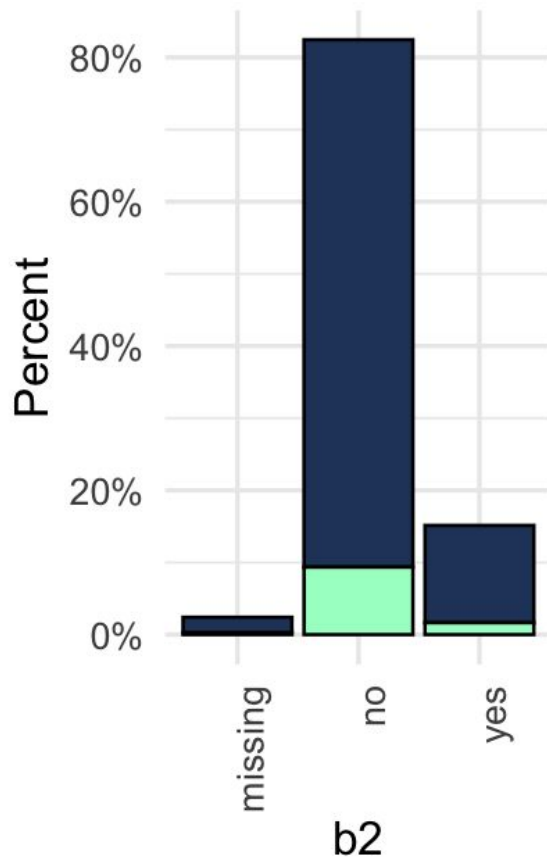


successful_sell no yes

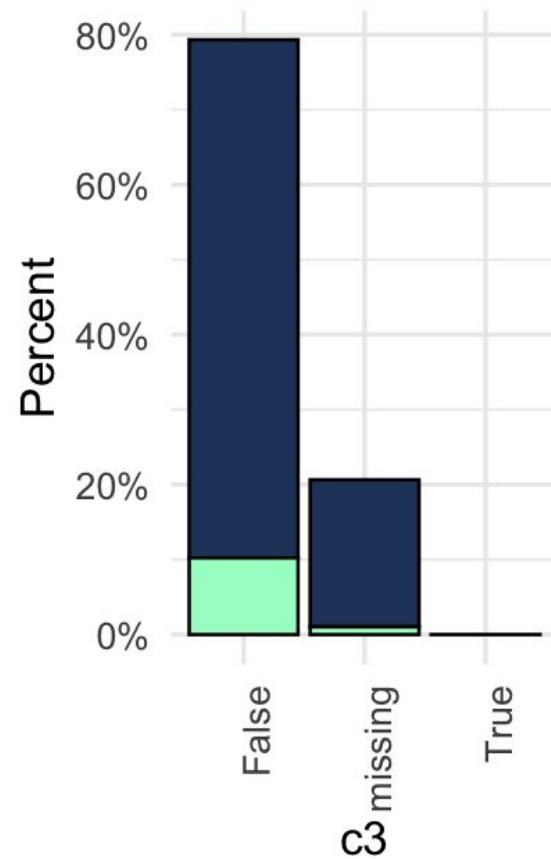
b1



b2



c3



successful_sell no yes