

Liya Zhang, Zhu Lyu

MA 346

Professor Carter

12/03/2020

## Final Project Documentation

For this project, our team is trying to predict if a client will subscribe to a deposit due to marketing effort. In this case, we are looking to predict the result of campaigns of a Portuguese banking institution (y), if the client will subscribe to a term deposit. The logistic regression model we came up with is composed of only a couple of variables documented by the bank: employee variation rate (emp.var.rate), last contact duration (duration), consumer price index (cons.price.idx), number of employees (nr.employed), contact communication type (contact), number of days that passed by after the client was last contacted from a previous campaign (pdays), consumer confidence index (cons.conf.idx), if the consumer has credit in default (default), number of contacts performed before this campaign for this client (previous). The process that we took is being published on <https://github.com/LiyaZhang-ziqing/Bank>. Below are steps we take to predict the variable 'y'.

The source of this bank marketing dataset is from UCI machine learning repository. We are using this dataset by splitting into training and validation datasets, so that training dataset is used to build the logistic regression model. And the validation dataset would be used to access our model, which would be further used in real-life if the model has a good performance. The model we build is looking to target marketing to the population that are most likely to accept deposit subscriptions in order to save resources from banks.

We took several steps in getting the data ready. After we load the dataset, we first drop NA fields inside the dataframe. Then, we dropped fields like jobs that are categorical values and cannot be turned into numerical or boolean values, because those are the acceptable forms for values in logistic regression modeling. Afterwards, we replaced the rest categorical values into boolean numeric or numerical values. The data is therefore ready to create the model.

Because we need to train and validate our model, we split the data into 60%(training) and 40%(validation). Secondly, we create a model using the training dataset using the scikit-learn's logistic regression tool. That is achieved by building a function, so that if the dataset changes, a new model can be built quickly using the same logic. We also build another function to test the model using the validation dataset. The values returned are the F1 score from the training dataset and validation dataset. F1 score represents how good our model is by using the precision and recall values. We can see that both F1 scores for training(0.4971) and validation(0.5019) datasets are relatively low. Possible reason is that we have a high number of predictor variables. It's likely that some predictors do not predict our response variable well.

Consequently, we try to create a better model and aim to achieve higher F1 scores. We first fit the model like we did in the last step. Then, we created a list of coefficients sorted in order of their importance to predict the model. With the list, we start using the most important variables and adding the rest one by one to see if the F1 scores would improve. After sets of trial and error, we came up with the best set of variables in predicting the variable y: 'emp.var.rate', 'duration', 'cons.price.idx', 'nr.employed', 'contact', 'pdays', 'cons.conf.idx', 'default', and 'previous'.

Lastly we decide to create a heatmap showing the visualization of the correlation between subscription status and its predictors. This visualization is being published on

<https://ancient-forest-16543.herokuapp.com/>. The part that we would like to focus on is the last column of the heatmap. It's showing that none of the predictors have strong correlation with the response variable, which potentially explain the relatively low F1 score from the logistic regression modeling. Moreover, one strong correlation piece that we notice is between employee variation rate and number of employees. We decided to keep both in our modeling process because the F1 score is about 2% higher than keeping just one of them.