

INTRODUCTION

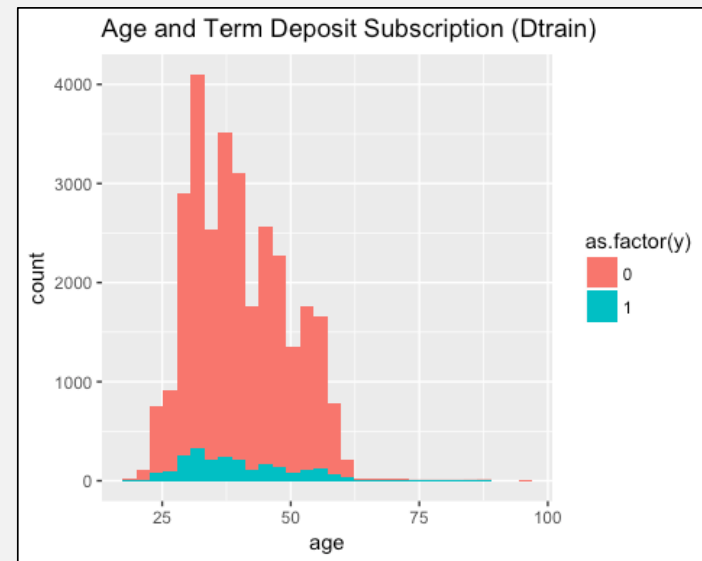
- Goal: predict the probability that a client would subscribe to a term deposit based on a variety of predictors
- Data:
 - Ten categorical and five integer variables
 - collected during a direct marketing campaign
 - 40000+ observations.

METHODOLOGY

- Subsetting: Data was divided into a training set and a validation set.
- Log-loss function used with the validation set to evaluate models.
- Approaches attempted:
 - XGB
 - LDA
 - Random forest
 - QDA
 - Clustering
 - Bruteforce to minimise log-loss
- **Linear model (logit LM chosen for final model)**

METHODOLOGY (CONT.)

- Data cleaning
- 'Age' and 'log(age)' were included to account for both non-linear and linear trends in the age data.
- Trial and Error
- Treatment of categorical variables:



RESULTS

- $$\log\left(\frac{\hat{y}}{1-\hat{y}}\right) = \frac{7.72}{(1.55)} + \frac{0.083age}{(0.014)} - \frac{3.45 \log(age)}{(0.576)} + \frac{0.14default_{no}}{(0.057)} -$$

$$\frac{0.25contact_telephone}{(0.098)} - \frac{1.28month_may}{(0.120)} - \frac{1.04month_jun}{(0.125)} - \frac{0.99month_jul}{(0.070)} -$$

$$\frac{1.20month_aug}{(0.078)} + \frac{2.24month_oct}{(0.275)} - \frac{1.09month_nov}{(0.084)} + \frac{1.137month_mar}{(0.134)} -$$

$$\frac{0.18day_of_week_mon}{(0.064)} + \frac{0.17day_of_week_wed}{(0.061)} + \frac{0.14day_of_week_thu}{(0.058)}$$
- All variables statistically significant
- AIC:15563

RESULTS

<u>EVALUATION METRIC</u>	<u>SCORE</u>
Log-Loss _(Validation)	0.255
False Negative Rate _(Valid.)	0.48
False Positive Rate _(Valid.)	0.077
Log-Loss _(Public Kaggle)	0.556
Log-Loss _(Private Kaggle)	0.580
Null Deviance	16511
Deviance	15533

Group: The Outliers

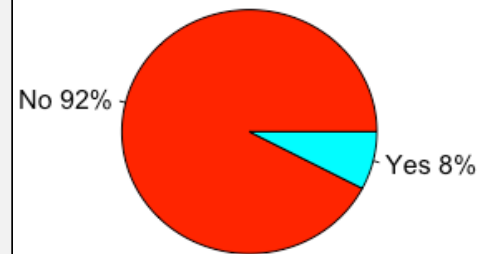
RESULTS - DISCUSSION

- Deviance < Null Deviance points towards model explaining the relationship within data
- False Negative Rate significantly greater than False Positive Rate
- Noting the months that appear to result in an increase or decrease in the probability of a client subscribing to a term deposit, is there seasonality?
- Clients contacted via cellular appear to have a higher probability of subscribing than those via telephone keeping all else constant
- Clients contacted in the middle of the week appear to have a higher probability of subscribing than those contacted earlier and/or later

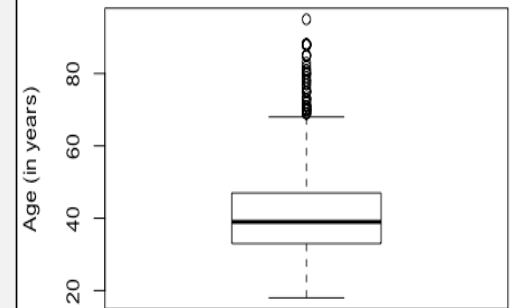
RESULTS – IN RETROSPECT...

- Possibility of overfitting?
- Class imbalance in variable 'y'
 - Solutions?
- Significant number of outliers in 'age' and 'campaign'
 - Solutions?
- In variable 'pdays', a value of 999 corresponded to someone not previously contacted.
 - This comprised the overwhelming majority of observations.
 - Solutions?

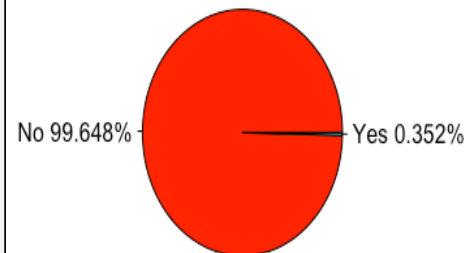
Pie Chart of Target Variable (Dtrain)



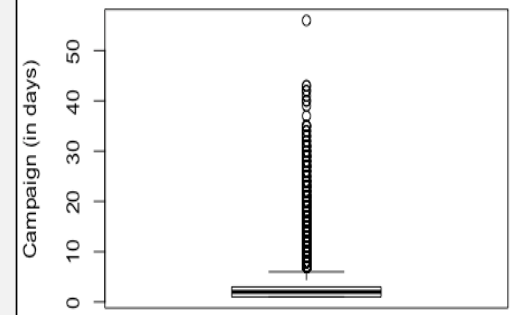
Boxplot of 'Age' (Dtrain)



Pie Chart of 'Not Previously Contacted' (Dtrain)



Boxplot of 'Campaign' (Dtrain)



Group: The Outliers

CONCLUSION

- A logit model comprising of significant predictors was produced
- This model could be used by businesses in order to identify the certain individuals that should be targeted through direct marketing
- Furthermore, this model could be used to identify certain strategy that could be implemented when direct marketing occurs
- Through analysis of the plots of various predictors, a variety of alternative methodology that perhaps should have been used were recognised for future models to consider