

Executive Report

Nicholas Lewis

5/5/2020

Introduction

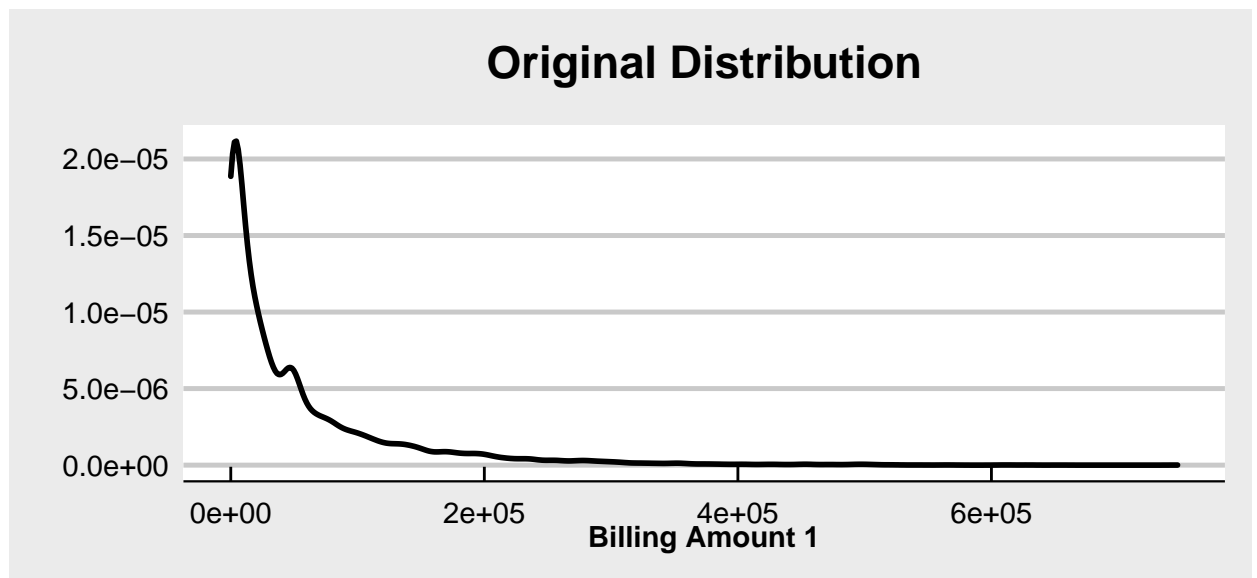
Data Processing

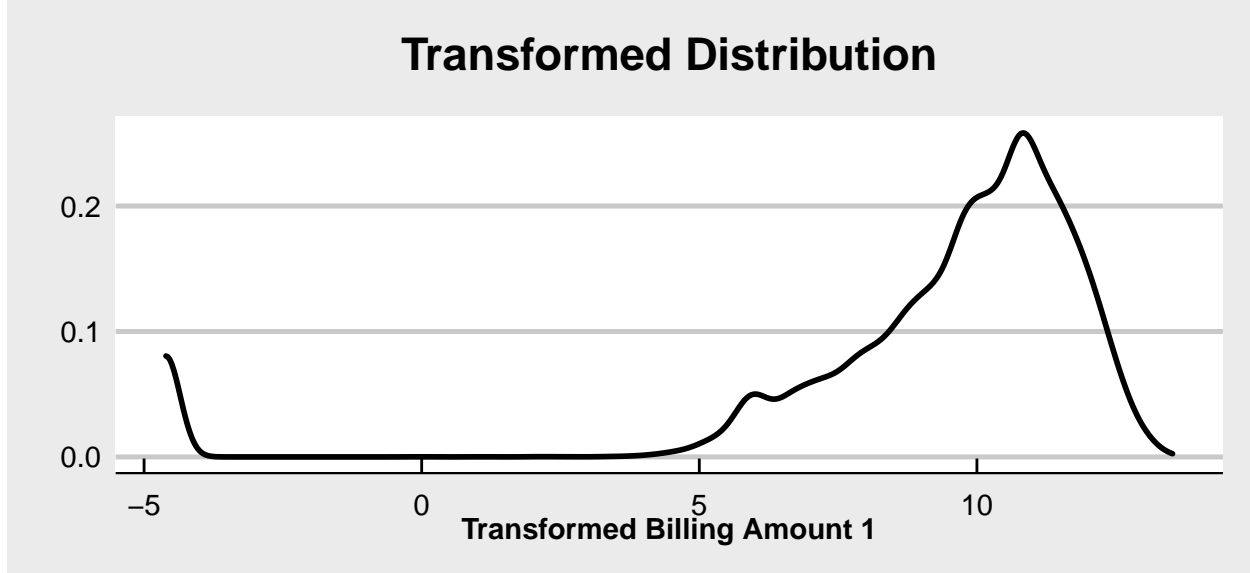
We first partitioned the data into two distinct sets: a training set and a testing set. The training data was used for all EDA and model building. Reported performance measures were calculated on the testing data.

Exploratory Data Analysis (EDA) and initial modelling attempts quickly illustrate the need for new features and transformation of existing features. First, simple density plots of the features ‘bill_amt1’, ..., ‘bill_amt6’, and ‘pay_amt1’, ..., ‘pay_amt6’ as well as age and balance limit are clearly not normally distributed. Logarithmic transformations are also taken in order to remedy this issue. There are several cases in which this transformation is problematic, however, since both the billing features and payment features contain zero values and bill amounts may also be negative. A slightly more complex version of this transformation is included below.

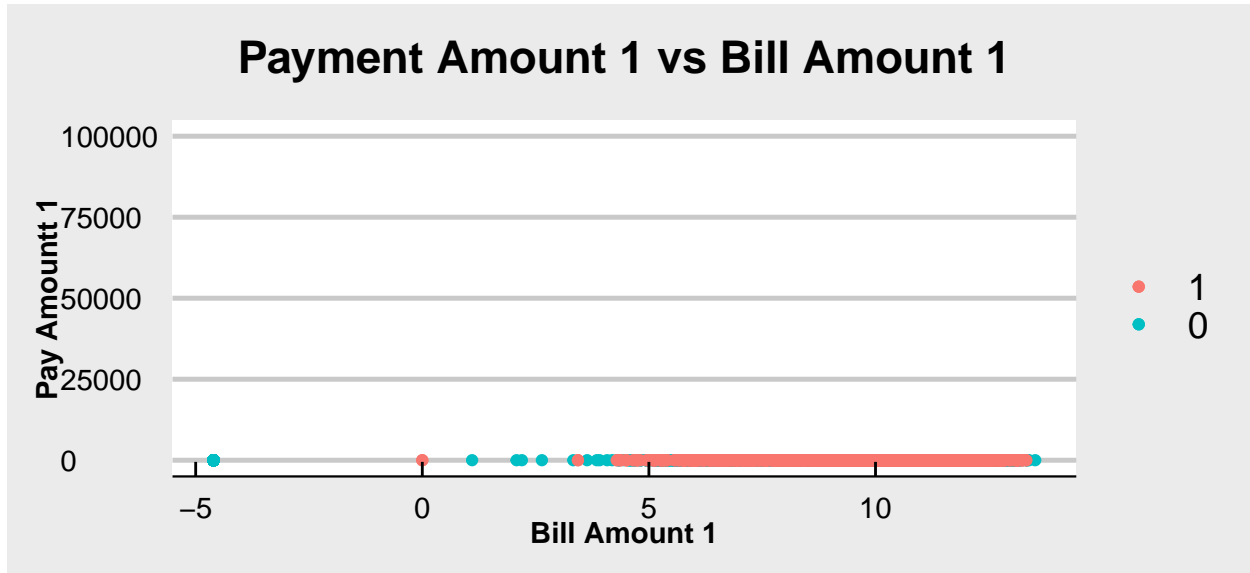
$$F(x) = \begin{cases} \log(x) & \text{if } x > 0 \\ \log .01 & \text{if } x = 0 \\ -\log|x| & \text{if } x < 0 \end{cases} \quad (1)$$

The transformation is illustrated graphically below for ‘bill_amt1’. All plots illustrating this transformation are included in Appendix -. It far from a perfect normal distribution even after the transformation, which may present issues with modelling techniques that assume normality.





Next, we create variables denoting the percent paid for each month. Intuitively, this ratio may be more influential than the billing and payment amount are separately. The plot below, as well as the rest of the family of plots included in Appendix B, illustrates that there is some sort of relationship between these two, but that observations who do not default do not adhere to this relationship very strongly. The noise extends above the top of this graph for those who do not default.



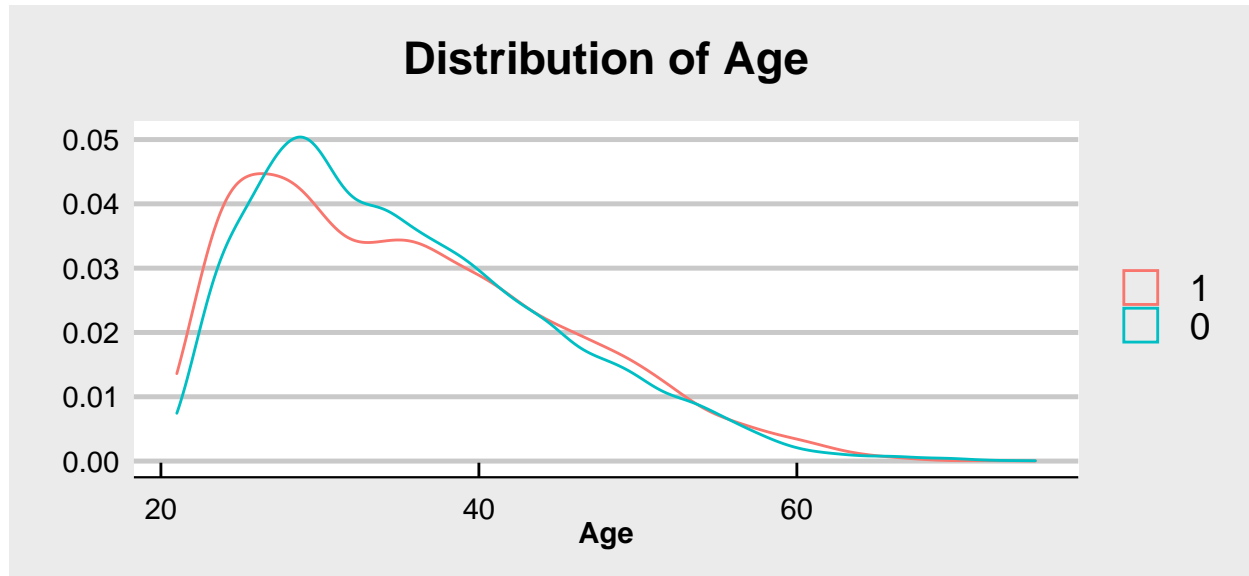
Again, observations which have a transformed billing amount of zero present a challenge. This is dealt with through the logic below.

$$F(x) = \begin{cases} \frac{\log(\text{Amount Paid})}{\log(\text{Amount Billed})} & \text{if } x = 0 \\ \frac{\log(\text{Amount Paid})}{\log(.01)} & \text{otherwise} \end{cases} \quad (2)$$

Extra features are also created for the log transformation of the mean billing amount over the six months and the log transformation of the mean payment amount. This is done to provide a more concise representation

of the billing and payment info if needed. The logic for the transformation of negative and zero amounts applies to this transformation as well.

Then, looking at the density plots of age by the value of the default value shows that there are age domains that are more likely to contain defaulting observations than others. To potentially reinforce these distinctions, we create indicator variables for the inclusion of an observation in the ranges $[0, 27)$, $[27, 40)$, $[40, 55)$, and $[55, \text{inf})$.



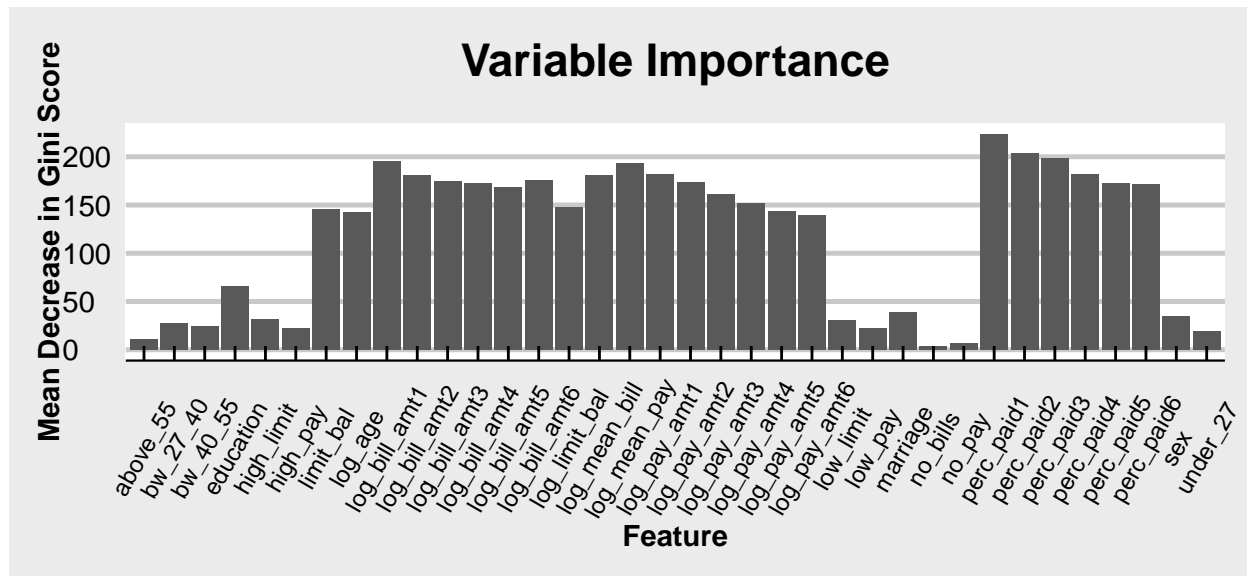
We use a similar methodology to create indicator variables for high and low pay (divided at \$2980) and for high and low billing limits (divided at $\$1.25 \times 10^5$). Similar plots for these variables are included in Appendix C.

Modeling

Random Forest

We first used the `train()` function supplied by the `caret` package to obtain the optimal values for the number of features considered at each split of each tree m and the number of trees. This process was fairly computationally expensive, but produced a final model containing 500 trees with $m = 2$.

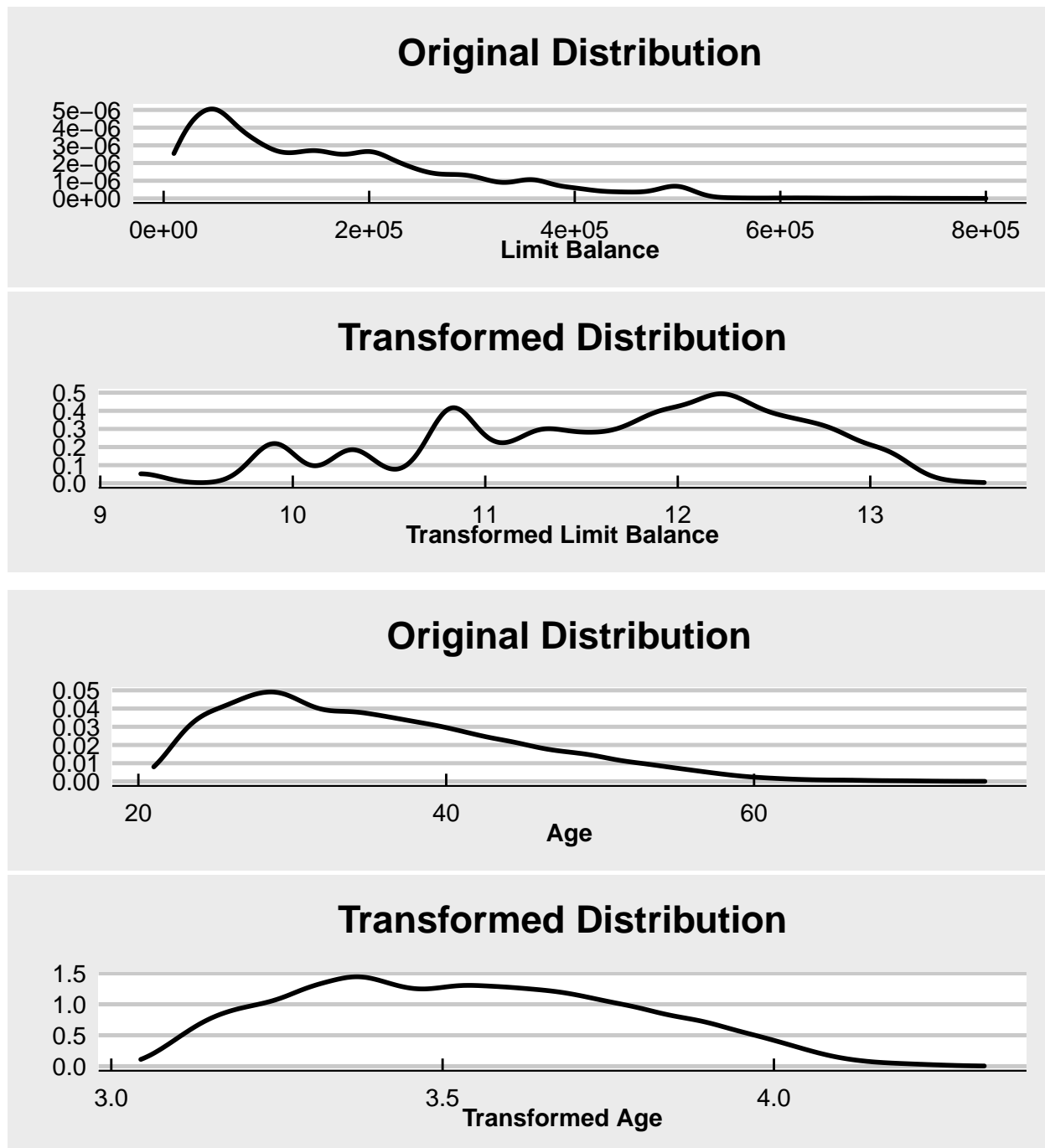
At this point, the log loss is fairly impressive, 0.4613 to be exact. However, we see that many of the variables created during the data processes are not important to the model (as displayed below). A low importance means that the feature was rarely used to make decisions in the forest because it was not a powerful predictor. However, because of the random nature, these low importance features will still be chosen at times. Removing them may allow more predictive features to be used more often, effectively increasing the power of our model.



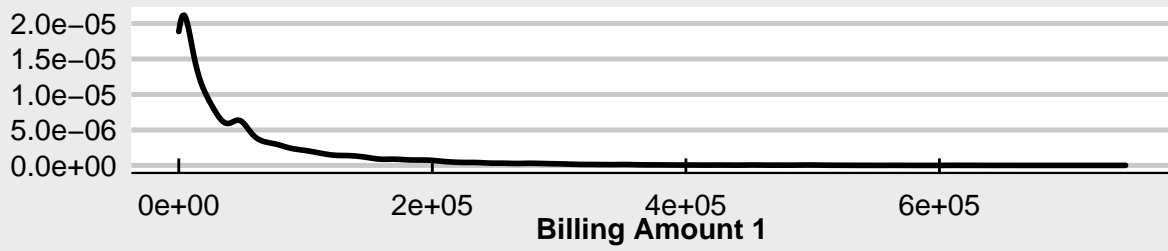
We remove the lowest variables, rebuild the model with $m = 2$ and $\# \text{ trees} = 500$. By doing this, we now obtain a log loss score of 0.4544.

Appendix A

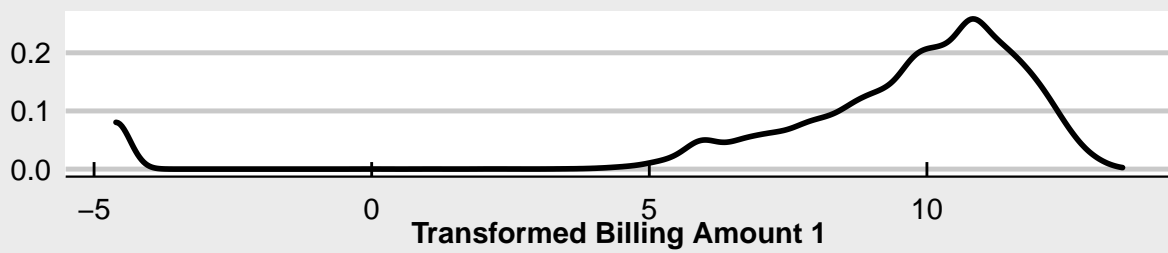
The following are plots of the variables transformed according to equation (1) before and after the transformation.



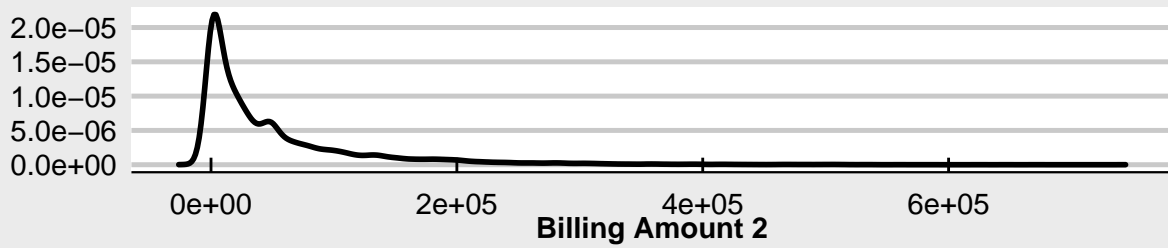
Original Distribution



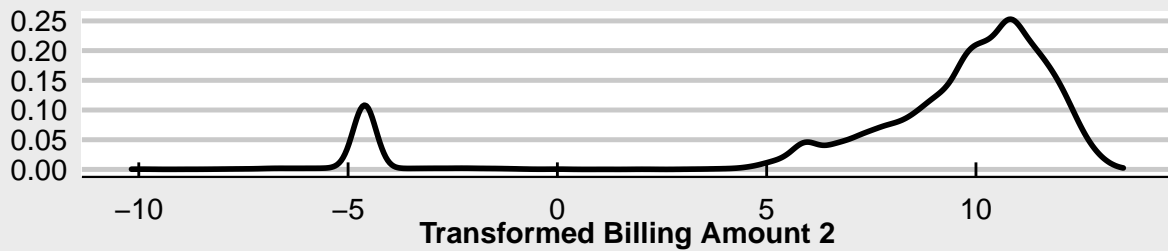
Transformed Distribution



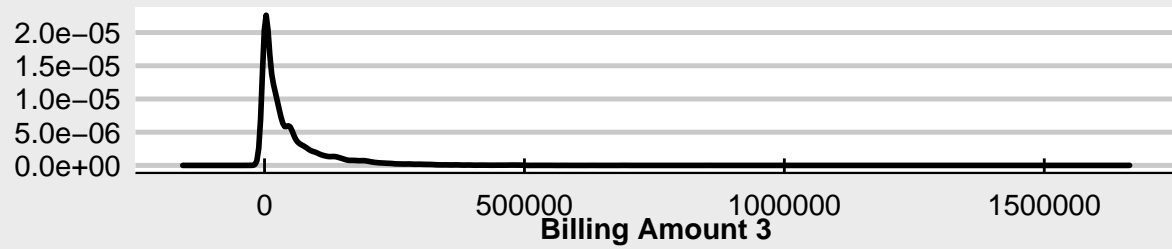
Original Distribution



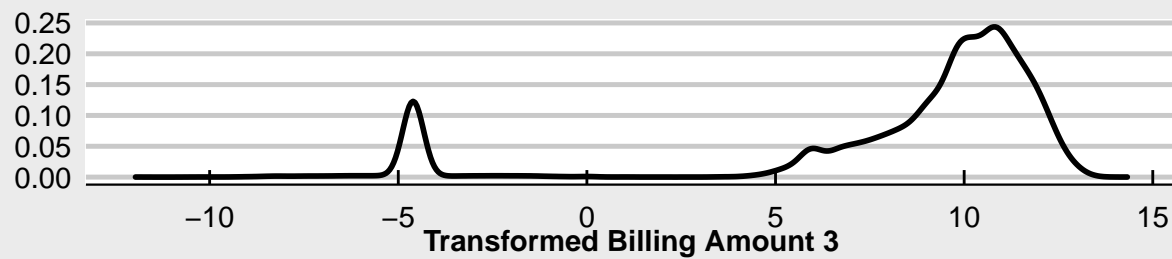
Transformed Distribution



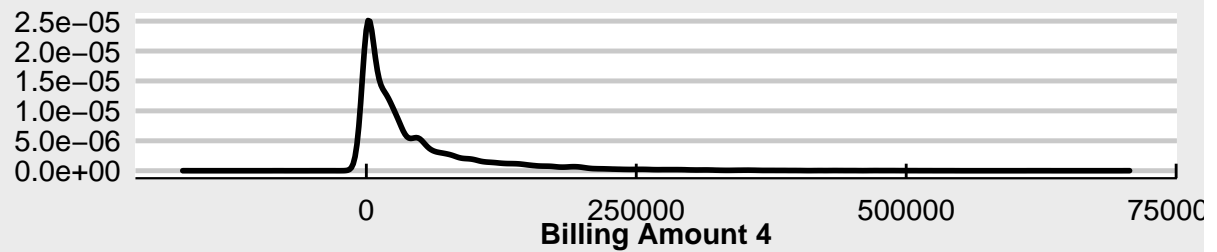
Original Distribution



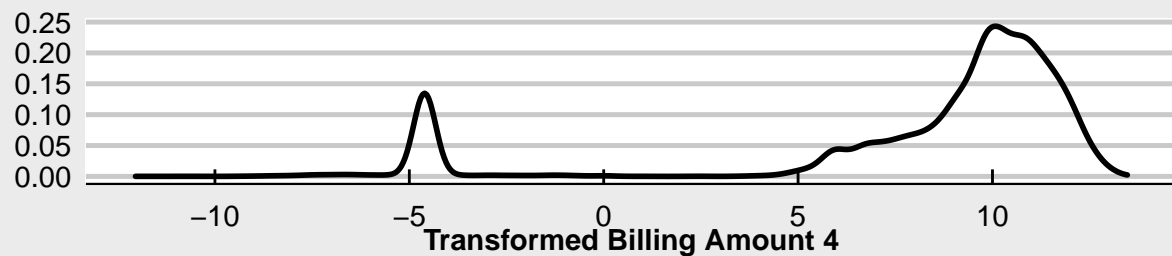
Transformed Distribution



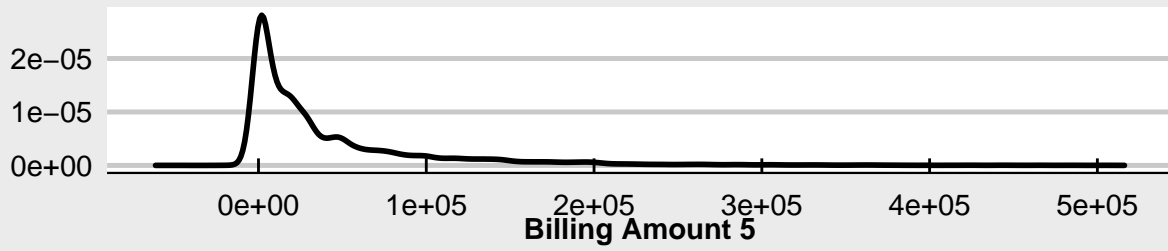
Original Distribution



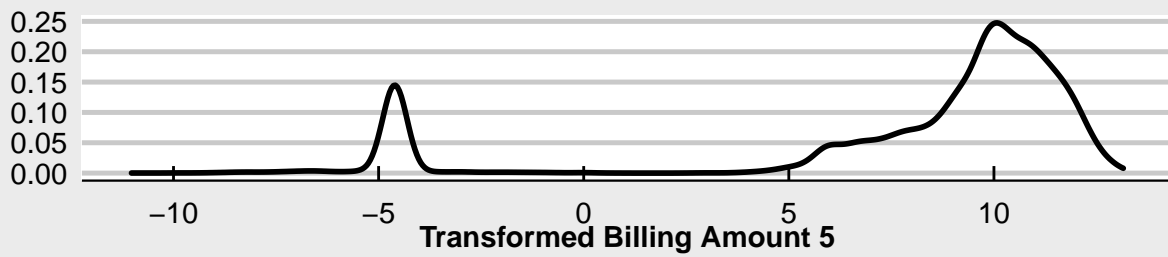
Transformed Distribution



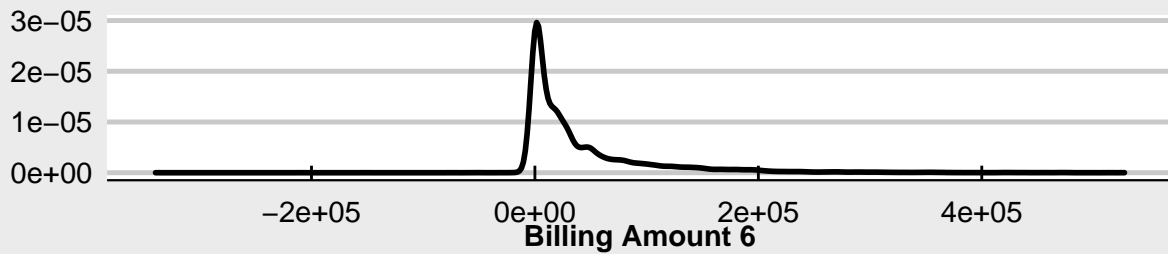
Original Distribution



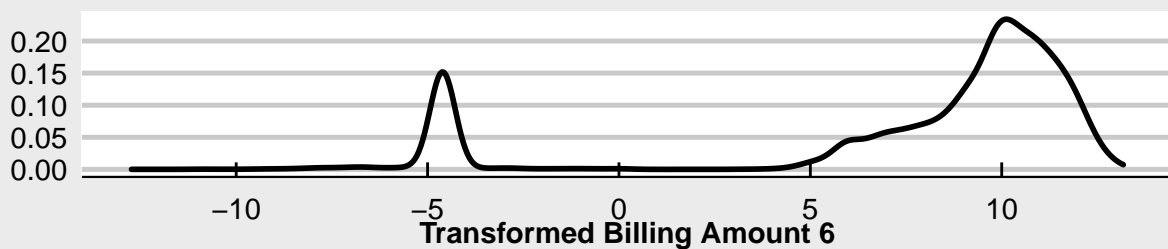
Transformed Distribution



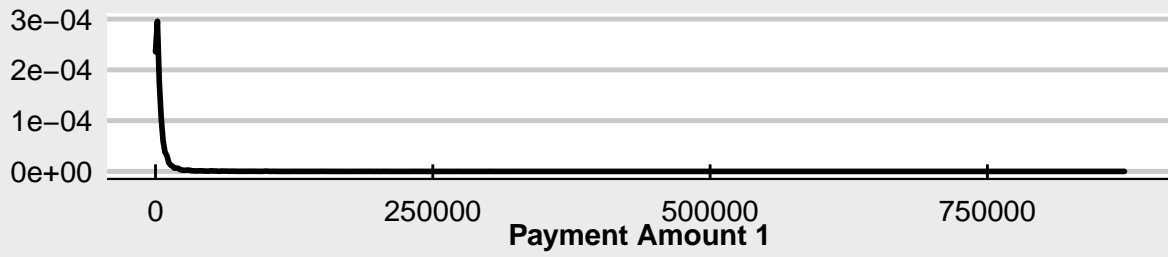
Original Distribution



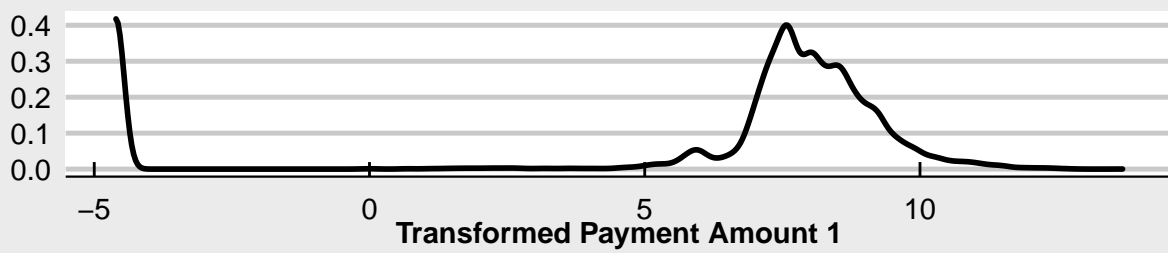
Transformed Distribution



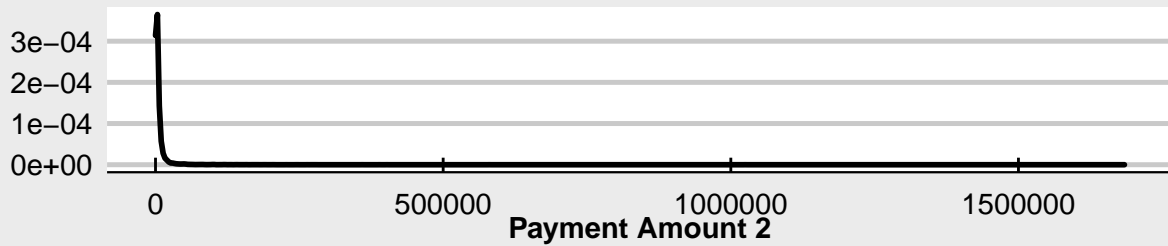
Original Distribution



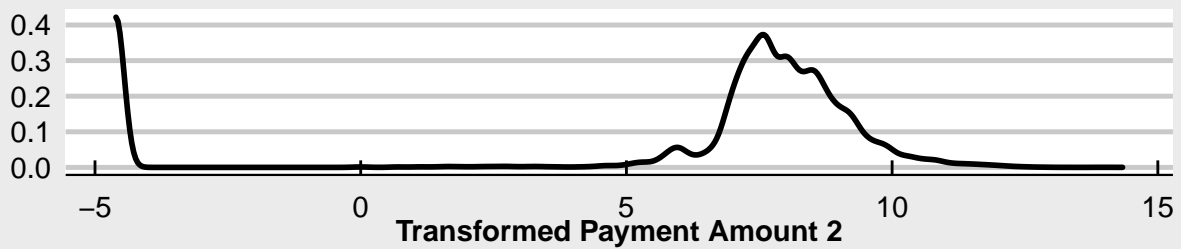
Transformed Distribution



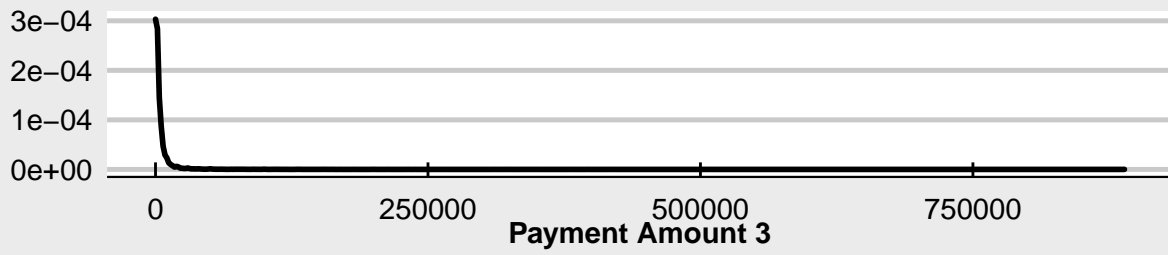
Original Distribution



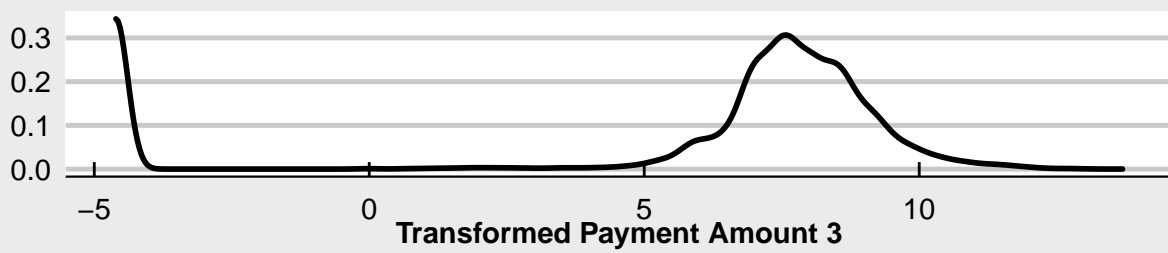
Transformed Distribution



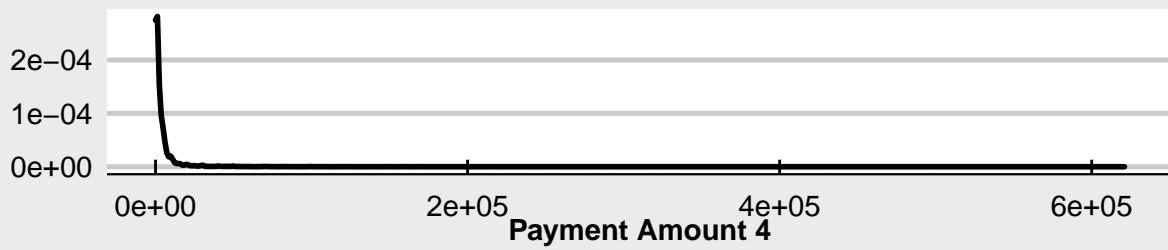
Original Distribution



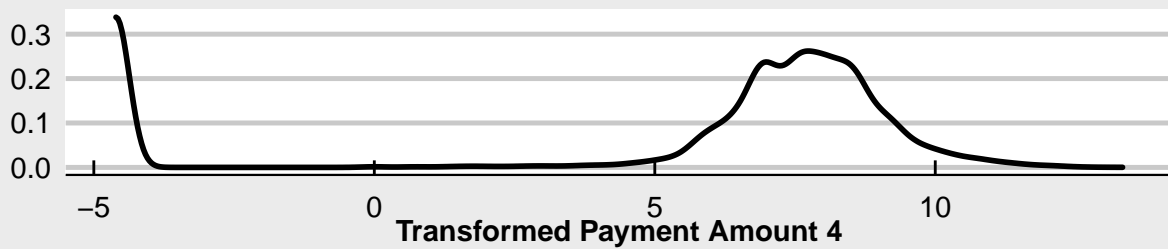
Transformed Distribution

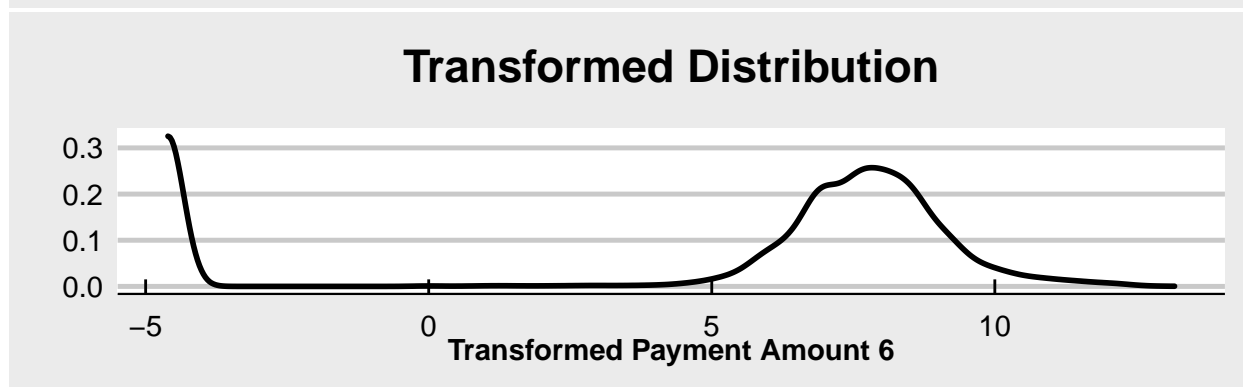
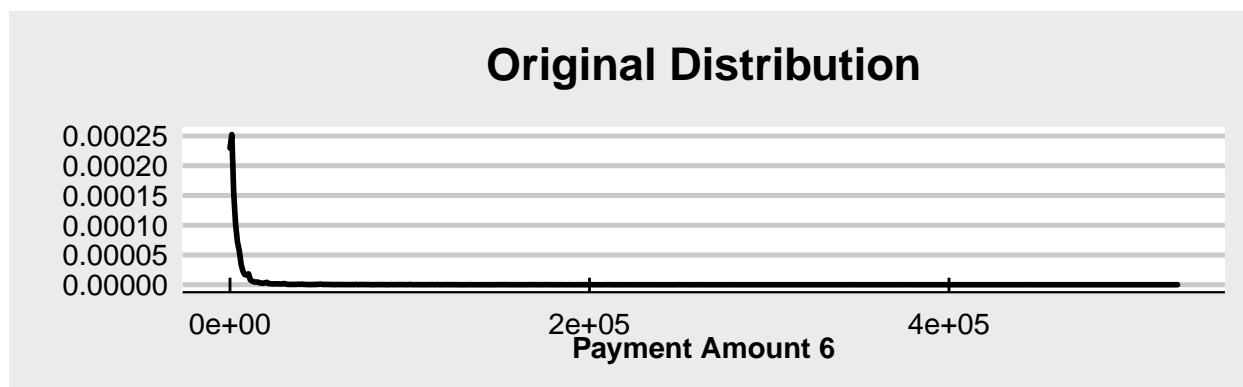
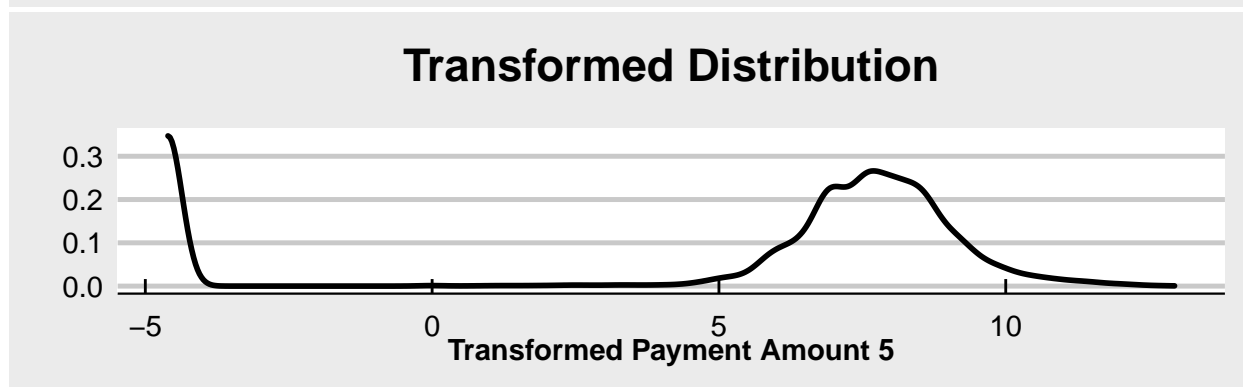
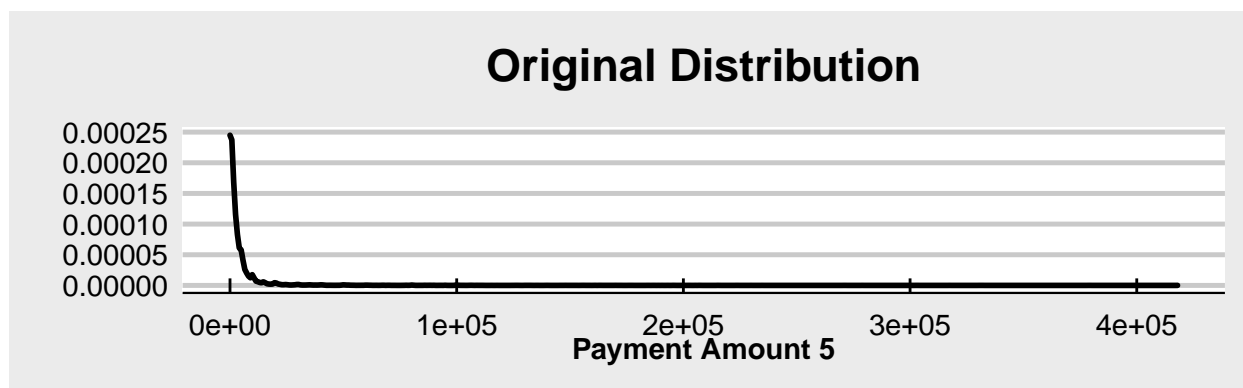


Original Distribution



Transformed Distribution





Appendix B

Will be the rest of the log_bill_amount vs log_pay_amount graphs