

Credit worthiness is a significant prediction problem for banks whose results can be improved through the use of Machine Learning techniques. Here, we analyze historical data for 150,000 customers provided by the client to train and develop a supervised logistic classifier. I go over the nuances below and achieve an accuracy of 93%. I then apply this to the test dataset of approximately 101,000 individuals.

A basic look through the dataset revealed substantial missing data: approximately 19.8% of entries were missing incomes and 2.6% of entries were missing dependents data. As these seem quite relevant to the prediction at hand, these entries were accounted for using mean and conditional mean methods so that we were able to use the full dataset, improving the predictive power of the other variables without harming our overall accuracy. By far the majority of individuals in the dataset make less than \$20,000 per month; the median is around \$5,400 for an annualized gross income of \$64,800. Additionally, after imputing data for child dependents, we see that a substantial number of individuals within the dataset have no children.

Indeed, looking through the data, we find that this seems to be a remarkably well-off group of individuals: no children, high median income, and an average of 8.4 credit lines or loans per customer. It is thus no surprise, then, that very few have been seriously delinquent within the past two years: only 6%.

Figure 1: Histogram of Monthly Incomes below \$20,000, after imputation

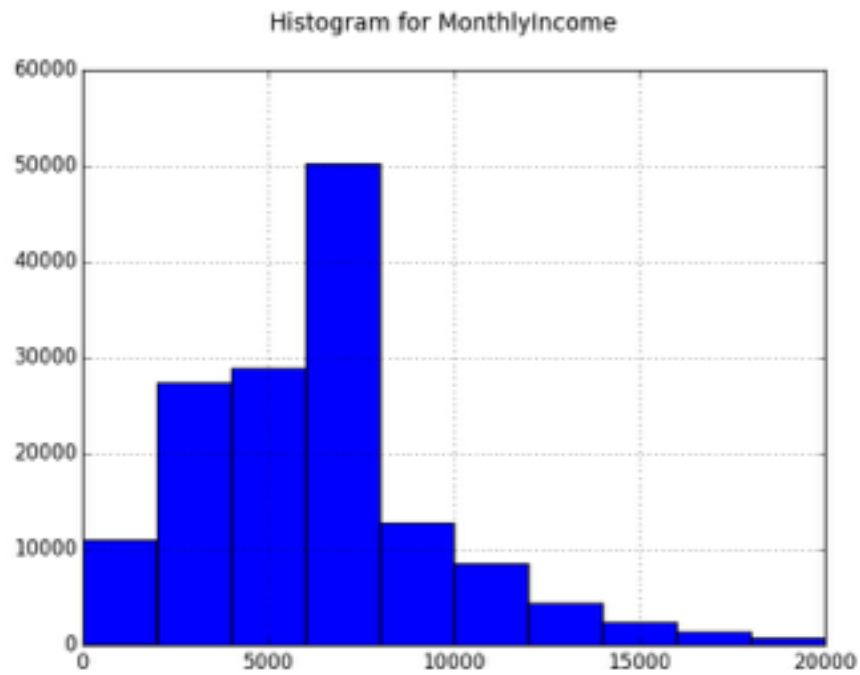
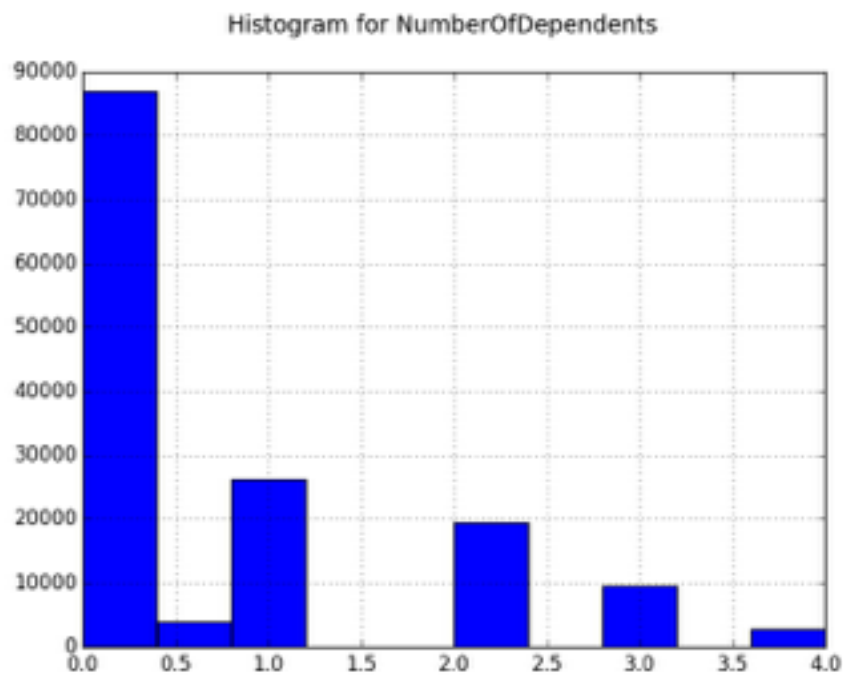


Figure 2: Histogram of Number of Dependents after imputation



Indeed, there is very little variation in most of the variables: Number of times past due on any measure of time, number of real estate loans or lines, and revolving utilization of unsecured lines. Even age was skewed older, with the mean age of 52—far above average.

For this reason, predicting new cases of default is expected to be difficult. I trained a model using the historical training set and achieved a 93.4% accuracy rate, which looks great. However, since 6.68% of individuals within the training set have been seriously delinquent, this means that guessing 'No serious delinquency' for all individuals within the dataset would result in a 93.3% accuracy rate. That is to say: the trained model only achieves a .1 percentage point increase in accuracy over always predicting no delinquency.

Looking at the coefficients, we also find that the data provided were not terribly helpful with our prediction. The two most impactful predictors were the number of times an individual was 30-59 days past due and number of times an individual was 90 days late. As mentioned before, neither of these variables has substantial variation, despite their predictive power. As such, we cannot expect reasonable accuracy in the model. Indeed, running on the test set, we see a predicted .004% rate of delinquency, which is too small to be taken seriously.

Next steps are to try another Machine Learning classification technique and see if improvements can be made. Additionally, provision of additional covariates around customers by the Client may be helpful; nevertheless, this may not be possible.

*Table 1: Coefficients from Logistic Model*

Variable	Coefficient
<b>RevolvingUtilizationOfUnsecuredLines</b>	-0.0000392
<b>age</b>	-0.030
<b>NumberOfTime30-59DaysPastDueNotWorse</b>	0.485
<b>DebtRatio</b>	-0.0000182
<b>MonthlyIncome</b>	-0.0000368
<b>NumberOfOpenCreditLinesAndLoans</b>	-0.008
<b>NumberOfTimes90DaysLate</b>	0.418
<b>NumberRealEstateLoansOrLines</b>	0.057
<b>NumberOfTime60-89DaysPastDueNotWorse</b>	-0.871
<b>NumberOfDependents</b>	0.086