



CREDIT SCORE MODELLING

ANDREW KINSMAN

MSc Data Analytics, Programming for Analytics 2015/16



Contents

Introduction	1
What are the Data?.....	1
Exploratory Data Analysis	1
Data Preparation.....	3
Decision Tree Model	4
Logistic Regression Model	5
Additional Models and Summary	7
References	7
Appendix 1 — Reflections.....	8
Appendix 2 — R Code	Error! Bookmark not defined.

Introduction

This project tackles a credit scoring problem that was the subject of an online competition run by Kaggle in 2011 [1]. Credit scoring is a standard procedure used by lenders to evaluate the potential risks of a loan, with a view to avoiding or reducing potential losses that would arise from bad debts [2]. Each client is given a credit score based on available information about that individual. Accurate credit scoring enables the lender to identify in advance which loans might be particularly high risk, so that they can perhaps act to prevent delinquency and/or make adequate provision for losses.

The task is to create a predictive model to determine whether or not an existing loan is a bad credit risk. Clearly there are only two possible outcomes: either the loan will become seriously delinquent (within a specified time period, in this case two years) or it won't. In this problem (known as a binary classification problem), the credit score of each loan represents the probability of it going bad. A credit score of close to 1 indicates a high credit risk, whereas a score close to zero is low risk.

The underlying idea of the model is to use data about loans for which the delinquency history is known to predict the credit risk of other loans, both current and future.

What are the Data?

Kaggle provide two sets of data: a training set for use in model creation and a test set which was used for assessing the competition entries. The training data contain 150,000 observations with 10 numerical explanatory variables or “predictors” (such as age, income, number of dependents, debt ratio). The test set has a further 101,503 observations with the same explanatory variables but without any values for the target variable, for which predictions need to be created. This target (or “dependent”) variable is whether or not the borrower experiences serious delinquency within the next two years. Kaggle also supply a data dictionary explaining each variable in the data [3].

Exploratory Data Analysis

Before starting any data preparation, it is advisable to undertake some exploratory data analysis (EDA), looking at the underlying structure of the dataset and performing summaries and visualisations. This will enable better understanding of the data and perhaps highlight some issues that may require special attention. Are there any missing or clearly incorrect values? How is each variable distributed (wide or narrow range, normally or skewed, with many outliers or not etc.)? Do the variables have the correct type (character, factor, integer, date)? Etc. Furthermore, EDA may reveal insights into possible relationships within the data and help identify which approach should be taken to the modelling process.

First the target variable is considered. Only around 10,000 (6.7%) of the 150,000 training set loans resulted in serious delinquency (*fig.1*). This severe imbalance between the two classes of the target variable could present problems when building a model, because the lack of minority class (“seriously delinquent”) data may

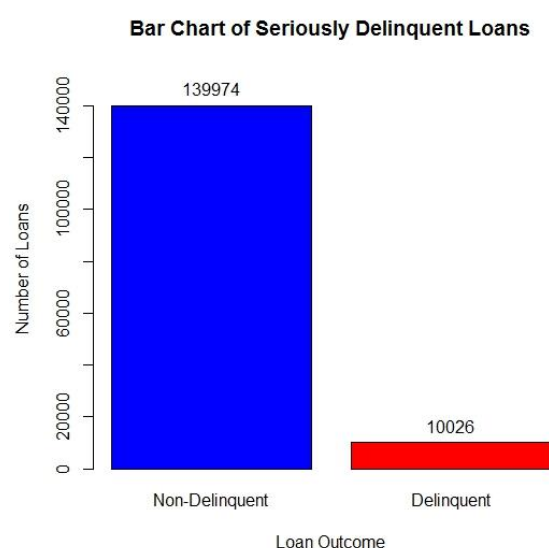


Figure 1

result in bias towards the majority classifier (“not seriously delinquent”). Such bias would cause over-optimistic estimates of the accuracy of the model, which would then fail to perform as well as expected when deployed on new data. Some methods for dealing with this imbalance will be discussed later in this report.

Further exploratory data analysis enables some initial observations to be made regarding the ten predictors, all of which (apart from one) are clearly skewed right:

1. The age variable has one erroneous value of zero, and might seem more or less normally distributed judging by the roughly bell-shaped histogram and fairly straight, diagonal middle part of the line in the Q-Q plot (*fig.2*). However, judging normality based on histograms is dangerous because any change in the number of bins will affect its shape. Indeed, a Lilliefors (Kolmogorov-Smirnov) test rejects the null hypothesis that age is normally distributed.

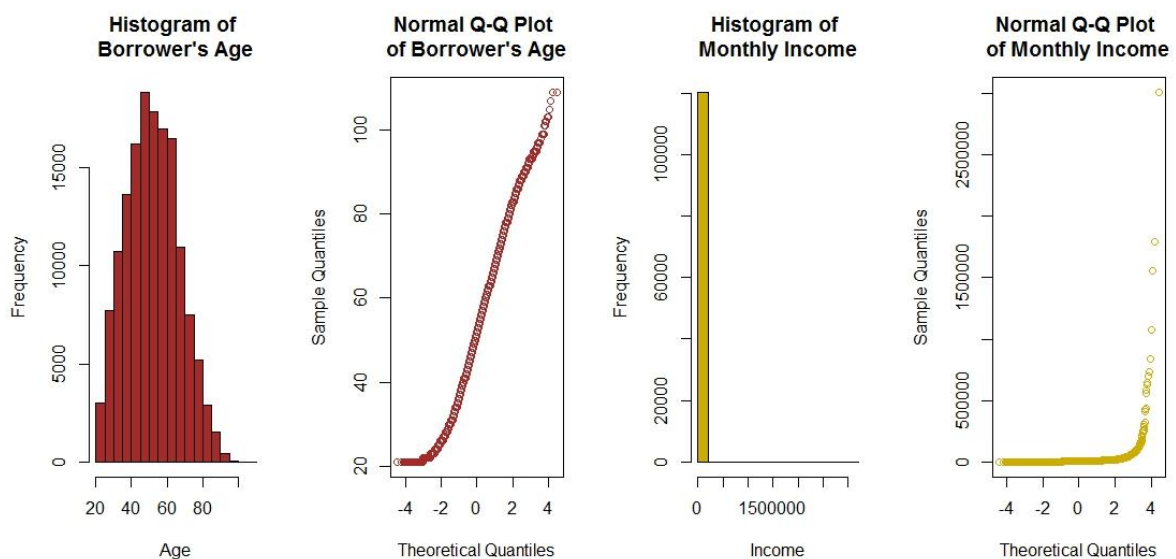


Figure 2

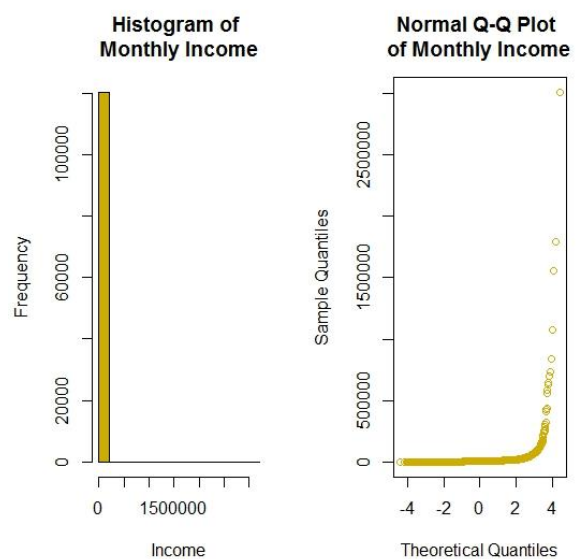


Figure 3

2. Unlike the age variable, monthly income has a heavily right-skewed histogram and no diagonal line in the Q-Q plot (*fig.3*). The 104 borrowers earning more than 1MM per annum (the currency is unspecified) help to squash the monthly income default R boxplot almost flat, so in order to produce a standard-looking boxplot the outliers must be removed (*fig. 4*).

The summary statistics reveal that 50% of clients earn between 3400 and 8249 per month (third quartile minus first quartile) and also that nearly 20% of the monthly income values are missing, which is clearly something that will require attention.

3. Most clients have between two and ten open credit lines and loans, but some have as many as 58 and there is clearly a right-skewed distribution.
4. The typical borrower has up to four dependents, but one has 13 and another has 20. There are nearly

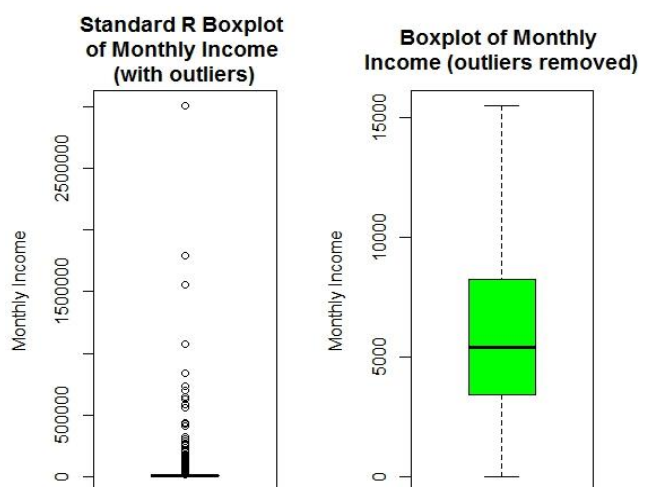


Figure 4

- 4,000 (2.6%) missing values that will need to be addressed.
5. Clients typically have either two, one or zero real estate loans or lines, although one outlier borrower has as many as 54.
 6. The third quartile of revolving utilisation of unsecured lines is just 0.56, but there are also some borrowers who possess hundreds or thousands of lines, going all the way up to over 50,000. This is one of several areas in which expert knowledge of credit scoring would be of assistance: How should these outliers be dealt with? Are they potentially important data points or simply erroneous data?
 7. The third quartile of debt ratio is just 0.9, but there are also some clients with a debt ratio into the thousands, going all the way up to over 300,000. Again expert interpretation would be helpful here.
 8. The other three variables, number of times 30-59 days, 60-89 days and 90+ days past due (late), each have the exact same count of outliers at values of 96 and 98. Both of these codes are commonly utilised in this kind of context to represent missing or corrupt data, coding issues etc.

Data Preparation

In this project the first data preparation step is to merge the training and test data together, so that dealing with any incorrect or missing values can be undertaken quickly and consistently on all data at the same time. There are various ways in which this cleaning can be carried out, including:

1. Rows with incorrect/missing values can be eliminated from the dataset.
2. Columns with incorrect/missing values can be removed from the dataset.
3. Some measurement of central location (e.g. the mean or median) can be substituted for the incorrect/missing values. Typically the median would be used as it is a more robust indicator of the central location than the mean, which is subject to influence by any extreme outliers.
4. A multiple imputation method can be used.

Rather than delete what might be useful data, an imputation approach was implemented. A simple approach of substituting the median was used for variables in which there were only a small amount of observations affected (replacing the “0” in the age column, and the “96”s and “98”s in each of the three columns relating to days late). However, while this approach is easy to implement, it may produce biased results in cases where the data is not missing completely at random. When there are larger numbers of missing or incorrect values (as with the monthly income and number of dependents variables), a multiple imputation approach is certainly preferable.

Multiple imputation works by generating predictions for each variable with missing values by utilising every other variable observation. It then uses those predictions to create plausible values for the missing data, with the process iterating until there is convergence over the missing values [4]. Several packages in R can perform multiple imputation, including the mice package that was used here to impute the missing values for monthly income and number of dependents.

Now that data preparation is complete, the data can be split back into the original training and test sets so that the model can be built on the training data and later assessed on the test data.

Having replaced the missing values, it is possible to construct a correlation matrix containing every variable in the training data. The strongest correlation (+0.43) is between the numbers of open credit and real estate loans or lines. The number of times past due for 30-59 day, 60-89 days and 90+ days are also fairly well correlated not only with the target variable, but also with each other. These

correlations between predictors indicate that there could possibly be a multicollinearity issue, whereby it would be problematic to distinguish between the individual effects of different variables.

Analysing these correlations, a boxplot of the number of times 30-59 days past due suggests that (not surprisingly) there is a much higher likelihood of delinquency for those who have been late with repayments in the past, whereas those who never fall behind on payments are extremely unlikely to become delinquent, as shown by the almost flat boxplot (fig.5).

The next step is to split the original training data in such a way that (say) 70% is randomly allocated to the train set on which the model will be built and the rest held back to form the validation set (which will subsequently be used to evaluate the models). Note that the split was executed in such a way that the same original proportions of serious delinquents and non-

delinquents are allocated to each set, since otherwise bias might be introduced simply because the individual train and validation sets are not adequately representative of the combined data.

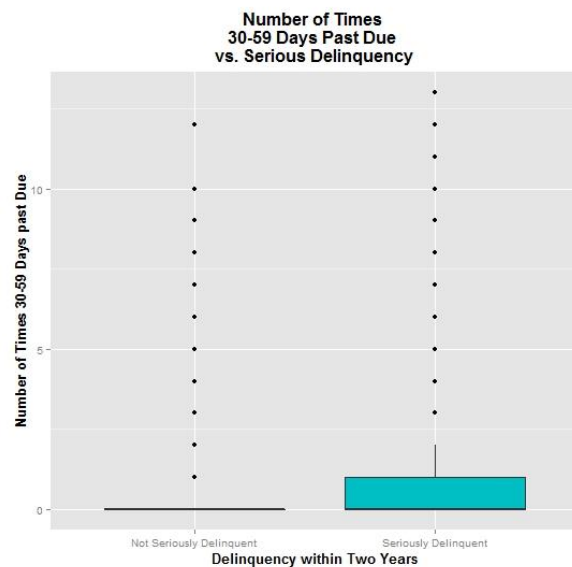


Figure 5

Decision Tree Model

Having cleaned the dataset, modelling can begin. There are many techniques for solving these dichotomous problems, and here a decision tree was tried first. Classification and regression tree (CART) models are popular because they are easy to interpret — predictions can be made for each new case simply by following a series of clear if-then rules. In the case of a classification tree, the tree algorithm works by iteratively splitting the dataset according to how homogenous the data is relative to the target variable [5].

In the tree produced using the rpart package (fig. 6), the most homogenous component was that 93% of loans in which at least one 90 or more days late were non-delinquent, and so forth until the model's stopping rules are met.

Starting at the top of the tree, the first root node (labelled 1), shows that 7% of loans ended in serious delinquency. The first split is then made by the decision tree, based on the number of times 90 days late variable. If this is greater than or equal to 0.5 (i.e. one or more, since this is an integer) then the serious delinquency rate is 41% (node 3), whereas otherwise it is just

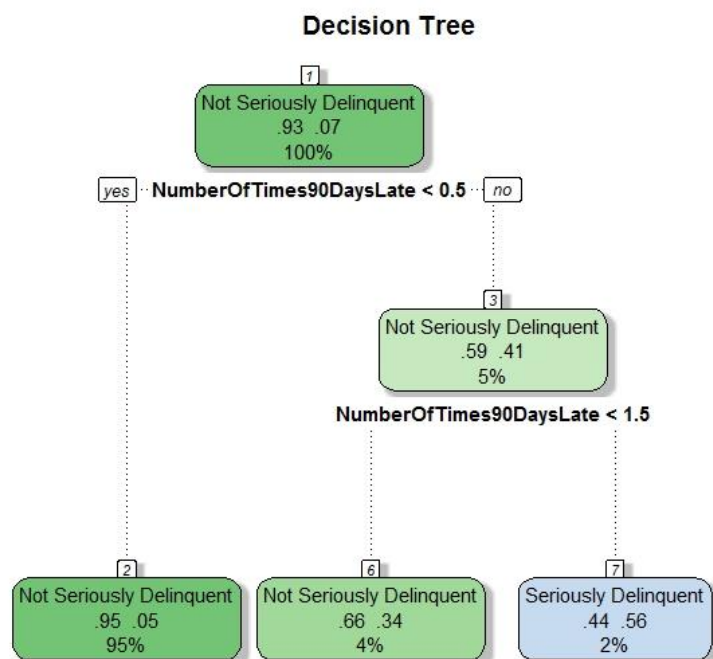


Figure 6

5% (node 2). The model makes no further splits on the left-hand side (node 2 is a "terminal" node).

Returning to the right-hand side of the tree, only one further split is made. Among those who have been 90 days late at least once, the ones who have been late more than once are more at risk, with a 56% serious delinquency rate compared to only 34% for those who have been that late only once.

This model suggests that just knowing whether or not someone has been 90 days late with a repayment either never, once, or more than once is a good indicator of whether they may become seriously delinquent. But how well does it make predictions? After running it against the validation data a classification (or "confusion" matrix) can be produced from the results (fig.

CONFUSION MATRIX FOR DECISION TREE			
		Predicted	
		Non-Delinquent	Delinquent
Actual	Non-Delinquent	41597	395
	Delinquent	2536	472

Figure 7

7). Model accuracy can be calculated by summing the correct predictions and dividing by the total number of observations. This gives 93.49% (42,069/45,000), slightly better than the baseline rate of 93.32% (which can be regarded here simply as the accuracy rate if you had simply guessed "non-delinquent" for every prediction).

Clearly there are some issues with this model. It lacks "precision", hitting on just 54.4% (472/ 867) predictions of serious delinquency, with a true positive rate (or "sensitivity") of just 472/3008 (15.7%). The model is poor at identifying actual serious delinquents, failing to provide the lender with accurate information on where their risk lies. In addition, the classification model yes/no approach doesn't allow any implementation flexibility — no probabilities are produced so (unlike in logistic regression, for example) there is no opportunity to adjust the "threshold" to allow for the fact that one kind of prediction error may be more costly than another.

Logistic Regression Model

Another common approach to binary classification is the logistic regression ("or logit") model, which predicts the *probability* of an event occurring depending on values of the explanatory variables. The regression coefficients for the predictors are calculated using maximum likelihood estimation, but unlike CART, logistic regression assumes a linear model, which may well not be appropriate here.

The first attempt at a logit model, trained on all of the explanatory variables (trainLog), runs into a complete separation issue [6]. The model does not converge because there are 18 clients with monthly income greater than 250,000, none of whom are seriously delinquent. The method chosen to resolve this was to split income into bins (five were chosen: below 1000, 1000-4999, 5000-9999, 10000-49999 and 50000 or more), treating it as a categorical rather than continuous variable.

In the revised model (trainLog1) every explanatory variable is significant at the 95% level except for revolving utilisation of unsecured lines and number of open credit lines and loans. It is worth considering whether one or more of those variables should be removed from the model, to reduce the chance of "overfitting" the training data, whereby the model fails to distinguish between the true "signal" of the underlying relationships and random noise.

One approach to model reduction is to remove the variable with the least significance and repeat the process until only "significant" variables are left. Another is to automate the process using stepwise regression based on some criterion. The Akaike Information Criterion (AIC), which builds a parsimonious model by penalising models that utilise more variables [7], was used here.

The AIC of the "full" model (trainLog1) is 41970. Can this be reduced with, say, a backwards stepwise regression model (trainLog2)? By removing the revolving utilisation of unsecured lines and number

of open credit lines and loans from the model, AIC can indeed be reduced to 41967. However, an ANOVA test on the two models produced a (barely) significant chi-square value (0.4658), suggesting that the second model does not fit as well (although for simplicity's sake it will be retained here).

It is also possible to perform forwards stepwise regression, starting with just the intercept, and this process may actually produce a different “best” model. Furthermore, stepwise regression does not consider every possible model, so the “optimal” model may even be overlooked [8].

Whereas the CART model produced an easy-to-understand decision tree, the coefficients from a logistic regression model are harder to interpret. To make a new prediction it is first necessary to input the values into the formula before converting the output into a probability. For example, let us consider a loan relating to a 40-year-old with one dependent and a monthly income of 4000 (with every other variable set to zero).

$$\begin{aligned}\text{prediction} &= -1.808513 + (40 \times -0.02768499) + (1 \times 0.04225236) + (0.03647074) \\ \text{prediction} &= -2.83719 \\ \text{odds} &= \exp(\text{prediction}) = 0.0585901 \\ \text{probability} &= \frac{\text{odds}}{\text{odds} + 1} = 0.0553473\end{aligned}$$

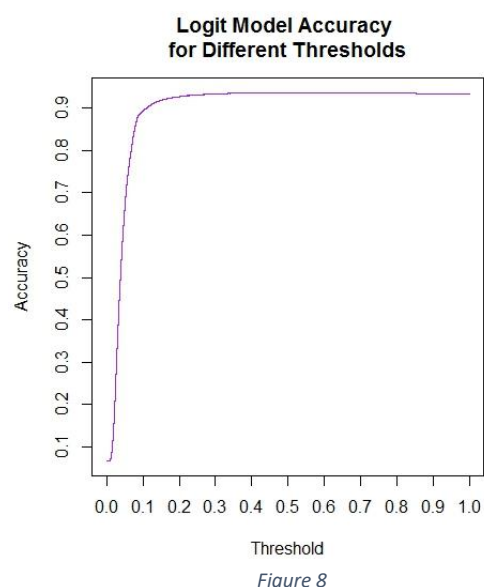
Thus there is roughly a 5.5% chance that the loan will become seriously delinquent.

By taking the exponent of the coefficients, it is possible to see how the odds of serious delinquency are affected by a change in one of the predictor variables (assuming that all others are held constant). Thus those odds are reduced by a factor of 0.973 for every one year increase in age, but increase by a factor of 2.361 for every time that someone has been more than 90 days late etc.

Compared to the decision tree, this model (assuming a “standard” 0.5 threshold) produces a slight improvement in accuracy, much better precision, and slightly worse sensitivity. However, that 0.5 threshold assumes that false positives (“Type 1” errors) and false negatives (“Type 2” errors) are equally costly. In the case of loans, this is not likely to be the case — false negatives (occasions when the model failed to predict serious delinquency that *did* occur) may be much more costly than false positives (when it predicted serious delinquency that *did not* occur, i.e. false alarms).

Given that (like the decision tree) this model underrepresents incidence of serious delinquency, it may be circumspect to lower the threshold so that more actual delinquency incidents are identified (higher sensitivity), even at the cost of worse specificity due to more false positives. Reducing the threshold to 0.2 here doubles the chance of correctly identifying actual delinquency with only a small trade-off in lost accuracy (*fig.8*). Ideally, however, the optimal threshold would be determined by someone with knowledge of the relative costs of these errors.

Instead of accuracy, which is not a robust metric for unbalanced data, the Kaggle competition used an absolute measure of the quality of predictions, the “area under the curve” (AUC). The AUC is derived from a Receiver Operating Characteristic (ROC) curve, which typically organises classifiers by plotting the true positive rate against the false positive rate. The ROC curve captures all thresholds simultaneously (and can thus assist with threshold selection) [9]. An AUC of 1 would be representative of perfect classification

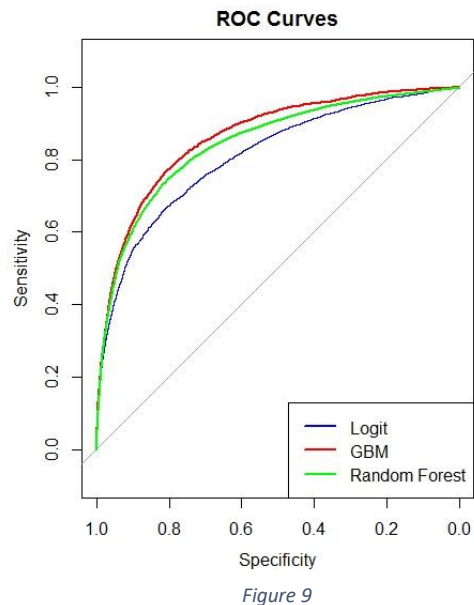


between high- and low-risk loans, whereas an AUC below 0.5 indicates that simply guessing would have been more effective. The competition winner achieved an AUC of 0.86956, but the decision tree and logit models here only managed 0.64829 and 0.81289 respectively on the exact same data.

Additional Models and Summary

In view of this poor performance, two alternative models were built: a gradient boosted machine (GBM) and a random forest. Each of these works by building ensembles of decision trees. Although more powerful than simple decision trees, they are much harder to interpret. Interestingly, however, the variable importance output of both new models suggests that revolving utilisation of unsecured lines (which was discarded using stepwise logit) is one of the key predictors. These models may perhaps be capturing some interactions that eluded the logit model.

Looking at the ROC curves (*fig.9*), the new models both perform well across a range of thresholds. The (blue) logit curve lies beneath the two other curves, indicating lower AUC. (N.B. The pROC package in R plots *specificity* on the x-axis rather than the false positive rate, i.e. $1 - \text{specificity}$, which means that the x-axis runs down from 1.0 in rather unintuitive fashion.



Here is a summary of the performance of all four models (*fig.10*):

TYPE OF MODEL	METRIC					
	Accuracy	Sensitivity	Specificity	AUC	Kaggle AUC*	Kaggle Position*
Classification Tree	0.93487	0.15691	0.99059	N/A	0.64829	868th
Logit (0.2 threshold)	0.92753	0.31848	0.97116	0.80964	0.81289	756th
GBM	0.93651	0.57197	0.94521	0.863	0.86413	417th
Random Forest	0.93558	0.56720	0.94234	0.8462	0.85335	635th

Figure 10 (Model Comparison)

* Kaggle results relate to the private leaderboard (for which submissions can still be assessed even though the competition has closed). This leaderboard contained 925 entries [10].

References

- [1] Kaggle, "Give Me Some Credit", available from <<https://www.kaggle.com/c/GiveMeSomeCredit>>.
- [2] "Credit score", (wiki article). Available from <https://en.wikipedia.org/wiki/Credit_score>.
- [3] Kaggle data available from <<https://www.kaggle.com/c/GiveMeSomeCredit/data>>.
- [4] Kabacoff, R. (2015) *R in Action (2nd edition)*, Manning, p 428.
- [5] Linoff G.S. & Berry, M.J.A. (2011) *Data Mining Techniques (3rd edition)*, Wiley p 237.
- [6] Albert A. & Anderson J.A. (1984) "On the existence of maximum likelihood estimates in logistic regression models", *Biometrika* 71: pp 1-10.
- [7] *Ibid.* 4, p 202.
- [8] *Ibid.* p 204.
- [9] *Ibid.* p 408.
- [10] Kaggle leaderboard from <<https://www.kaggle.com/c/GiveMeSomeCredit/leaderboard>>.

Appendix — Reflections

This project presented some interesting challenges. Initial analysis revealed some missing values and coding issues, and also that there was a severe imbalance within the target variable, which could result in model bias towards the majority class.

In order to predict which loans were higher risk, binary classification modelling was used. Credit scores were generated for each loan, based on the probability of it becoming seriously delinquent within two years.

After preparing the data, various different models were tried with mixed results. The classification decision tree was easy to interpret, but was poor at identifying actual incidence of serious delinquency and failed to generalise well when applied to new data. Although the logit model achieved better AUC performance, it lacked the intuitive interpretability of the decision tree. (Both of these models also underrepresented the actual quantity of serious delinquency incidence in their predictions, indicating class bias.) The GBM and random forest models were even more difficult to interpret, since they are essentially “black box” techniques, but each produced far superior performance, which could possibly even be enhanced by some parameter tuning.

Access to domain expertise would also be extremely useful. A credit scoring expert would be able to fully dissect the structure found in each variable, identify which outliers might be erroneous data, and advise ways in which predictors might be combined to generate useful derived variables. In order to assess the potential value of this, a new variable was generated by dividing income by family size (the logic being that larger families might possibly have relatively less discretionary income), which resulted in slightly better GBM performance, an AUC of 0.86445 on the Kaggle competition test set.

More sophisticated models would incorporate some mechanism for offsetting target variable imbalance, perhaps by oversampling the minority class (or possibly undersampling the majority class, although that would involve throwing out data) to balance the dataset. This oversampling would typically involve bootstrapping the minority class to generate many more observations sourced from the existing ones (bootstrapping is a form of random sampling with replacement). Another technique for dealing with unbalanced data is to adjust the costs of misclassification, so that the penalty for missing actual incidence of delinquent loans is higher, but changing the threshold becomes academic if the goal is simply to maximise AUC (as was the case in this competition).

Another common technique for potentially enhancing performance is to create an ensemble from multiple models, combining their predictions using an averaging process. Models often produce better predictions when combined than they do individually (indeed, data competitions are frequently won by people or teams using ensembles). A simple average of the predictions from the GBM and random forest models from this project resulted in an improvement in AUC to 0.86534, good for 221st place and within .005 of the winner’s score. Naturally, in the financial world improving AUC by just a fraction might be worth tens of thousands of pounds.

With regards other applications of these models, none of the other project proposals from the course clearly lend themselves to binary classification analysis. However, these modelling techniques are common when the target variable is dichotomous, in areas of medicine (e.g. predicting whether a tumour is malignant or benign), customer relationship management (e.g. predicting whether or not someone will renew a mobile phone contract), polling (e.g. predicting whether someone will vote Republican or Democrat), parole decisions (predicting whether or not someone will violate their parole) etc.