

CLASSIFICATION BASED CREDIT RISK ANALYSIS: THE CASE OF LENDING CLUB

AADI GUPTA *

PRIYA GULATI †

SIDDHARTHA P. CHAKRABARTY ‡

Abstract

In this paper, we perform a credit risk analysis, on the data of past loan applicants of a company named Lending Club. The calculation required the use of exploratory data analysis and machine learning classification algorithms, namely, Logistic Regression and Random Forest Algorithm. We further used the calculated probability of default to design a credit derivative based on the idea of a Credit Default Swap, to hedge against an event of default. The results on the test set are presented using various performance measures.

Keywords: Credit risk; Classification algorithm; Exploratory data analysis

1 INTRODUCTION

Lending Club, headquartered in San Francisco, was the first peer-to-peer lending institution to offer its securities through the Securities and Exchange Commission (SEC) and enter the secondary market [1]. This article is essentially a case study, on how financial engineering problems can be addressed using Machine Learning (ML) and Exploratory Data Analysis (EDA) approaches. Lending Club specializes in extending different types of loans to urban customers, which is decided on the basis of the applicant's profile. Accordingly, the data considered in this work, contains information about past loan applications and whether they were "defaulted" or "not". One of the main objectives of this work is the calculation of the Probability of Default (PD) and (in case of a default) the determination of Loss Given Default (LGD), Exposure at Default (EAD), and finally, the Expected Loss (EL), making use of the historical data. Also, using the "Recovery Rate" (estimated while calculating "LGD" from "EAD") and the PD, a simple credit derivative is implemented, based on the concept of Credit Default Swaps (CDS) which can be used to hedge against such defaults. It may be noted that the total value a lender is exposed to when a loan defaults, is the EAD and the consequent unrecoverable amount for the lender is the LGD [2]. Accordingly, EL is defined as,

$$EL = EAD \times LGD \times PD.$$

The driver of this work is the notion of "Classification Algorithm", which weighs the input data (of the applicant, in this case) to classify the input features into positive and negative classes [3] (default on the loan or not, in this case). Accordingly, we consider two "Classification Algorithms", namely, the Logistic Regression and the Random Forest.

The main idea behind Logistic Regression is to determine the probability of a particular data point belonging to the positive class (in the case of binary classification) [4]. The model does so by establishing a "linear relationship" between the independent and the dependent variables. The weights of the linear relations is determined through

*Department of Mathematics, Indian Institute of Technology Guwahati, Guwahati-781039, India, e-mail: aadi18@alumni.iitg.ac.in

†Department of Mathematics, Indian Institute of Technology Guwahati, Guwahati-781039, India, e-mail: gulati18@alumni.iitg.ac.in

‡Department of Mathematics, Indian Institute of Technology Guwahati, Guwahati-781039, India, e-mail: pratim@iitg.ac.in

the minimization of a Loss Function, which is achieved by using the “Gradient Descent Optimization Algorithm”. Since Logistic Regression first predicts the probability of belonging to the positive class, therefore it creates a linear decision boundary (based on a *threshold*, set by the user), separating the two classes from one another. This decision boundary can now be represented as a conditional probability [5]. Implementation of the Random Forest algorithm involves the training stage construction of several decision trees [6], and predictions emanating from these trees are averaged to arrive at a final prediction. Since the algorithm uses an average of results to make the final prediction, the Random Forest algorithm is referred to as an ensemble technique. Decision Trees are designed to optimally split the considered dataset into smaller and smaller subsets, in order to predict the value being targeted [7, 8]. Some of the criterion used to calculate the purity or impurity of a node, include Entropy and Gini Impurity. In summary, a decision tree splits the nodes on all the attributes present in the data, and then chooses the split with the most Information Gain, with the Decision Tree model of Classification and Regression Trees (CART), being used in this paper.

The firm based model uses the value of a firm to represent the event of default, with the default event being represented by the boundary conditions of the process and the dynamics of the firm value [9]. In particular, we refer to two well-established models. Firstly, we mention the Merton model based on the the seminal paper of Black and Scholes [10], which is used to calculate the default probability of a reference entity. In the context, the joint density function for the first hitting times is determined [11]. Secondly, we have the Black-Cox model, which addresses some of the disadvantages of the Merton model [10]. In order to hedge against credit risk, the usage of credit derivatives is a customary approach [12, 13, 14], with Credit Default Swaps (CDS) being the most common choice of credit derivatives. This type of contracts entail the buyer of the CDS to transfer the credit risk of a reference entity (“Loans” in this case constitute the reference entity) to the seller of the protection, until the credit has been settled. In return, the protection buyer pays premiums (predetermined payments) to the protection seller, which continues until the maturity of the CDS or a default, whichever is earlier [10]. The interested reader may refer to [10] for the formula of CDS spread per annum, to be used later in this paper. Another widely used credit derivative, albeit more sophisticated than the CDS are the Collateralized Debt Obligation (CDO), which is a structured product, based on tranches [14]. CDOs can further classified into cash, synthetic and hybrid. The interested reader may refer to [14] for a detailed presentation on pricing of synthetic CDOs.

2 METHODOLOGY

The goal of this exercise is the approximation of a classification model, on the data considered (from the peer-to-peer lending company Lending Club), in order to predict as to whether an applicant (whose details are contained in the considered database) is likely to default on the loan or not. Accordingly, to this end, it is necessary to identify and understand the essential variables, and take into account the summary statistics, in conjunction with data visualization. The dataset used for the estimation of PD, EAD, LGD and EL was obtained from Kaggle [15]. The data contains details of all the applicants who had applied for a loan at Lending Club. There were separate files obtained for the accepted and the rejected loans. The file “accepted loans” was only used, since the observations were made on the applicants who ultimately paid the loan, and those who defaulted on the loan. The data involved the details of the applicants for the loan, such as FICO score, loan amount, interest rate, purpose of loan etc. Machine Learning (ML) algorithms were applied on this data to predict the PD after data exploration, data pre-processing and feature cleaning.

The feature selection was performed on the considered data, bearing in mind the goals of predictive modelling.

Given the large number of existent features, it was desirable to achieve a reduction in their numbers, thereby retaining only the relevant ones. This action abetted the reduction in computational cost and improvement in the performance of the model.

The data pre-processing was carried out before the application of the classification model, in order to obtain the probabilities of default. This exercise involved the removal or filling of missing data, removing features which are repetitive or deemed unnecessary and the conversion of categorical features to dummy variables. Some of the dropped columns were “emp_title”, “emp_length”, “grade”, “issue_d” and “title”. After pre-processing, the data was divided into the customary training set and test set. The former was then used to train the classification model, with the results being evaluated subsequently, on the latter.

In order to accomplish the goal of comparison of different models, the problem of predicting the PD (denoted by p) was treated as a classification problem, where, if $p \geq 0.5$, it is treated as 1 (*i.e.*, the person will default) and if $p < 0.5$, it is treated as 0 (*i.e.*, the person will not default). Accordingly, $p = 0.5$ was used as the *threshold* for the classifier.

It was observed from the data, that the classes are imbalanced, and therefore, accuracy is not a good metric, when it comes to comparison of different models. Accordingly, we instead used the “Confusion Matrix” (to evaluate the performance of our model), which is a table with four different combinations of actual and predicted values *i.e.*, True Positive (TP), True Negative (TN), False Positive (FP) and False Negative (FN). It is beneficial for measuring *Recall* $\left(\frac{TP}{TP+FN}\right)$, *Precision* $\left(\frac{TP}{TP+FP}\right)$ and *F1-score* (Harmonic Mean of Recall and Precision).

Receiver Operating Characteristic (ROC) curve is used to compare the performance of a classification algorithm by varying the threshold value [16]. ROC is a probability curve and the area under the curve (AUC) represents the likelihood (Random Positive has dominant ranking over Random Negative) *i.e.*, the degree of separability. The AUC represents how successful is the model in differentiating between the classes. The higher the AUC value, the more amenable is the model at classifying 0 classes as 0, and 1 classes as 1 (*i.e.*, correct prediction). The ROC curve was plotted with TPR (True Positive Rate) against the FPR (False Positive Rate) where FPR is on the x -axis and TPR is on the y -axis. Note that TPR is the same as *Sensitivity* and $1-FPR = \frac{TN}{FP+TN}$ is also called *Specificity*.

3 RESULTS

The statistics-based feature selection methods were used to assess the correlation between every input variable and the target variable, and those input variables which demonstrated the strongest relationship vis-a-vis the target variable, were selected. Further, for selecting the relevant input variables, a comparative analysis was performed and some of the observations, as a result of this exercise, are noted below:

- (A) It looks like F and G sub-grades do not get paid back that often.
- (B) The loans with the highest interest rates are faced with a greater likelihood of default
- (C) Only 339 borrowers have reported an annual income exceeding 1 million.
- (D) It seems that the more the “dti”, the more likely is it that the loan will not be paid.
- (E) Only 1211 borrowers have more than 40 open credit lines.
- (F) Only 1299 borrowers have more than 80 credit lines in the borrower credit file.

The graph in Figure 1 shows the correlation between various features, like int_rate and loan_status (*i.e.*, whether the debtor will default or not). Positive correlation implies that more the value of particular feature, more is the chance that the debtor will default on their loan.

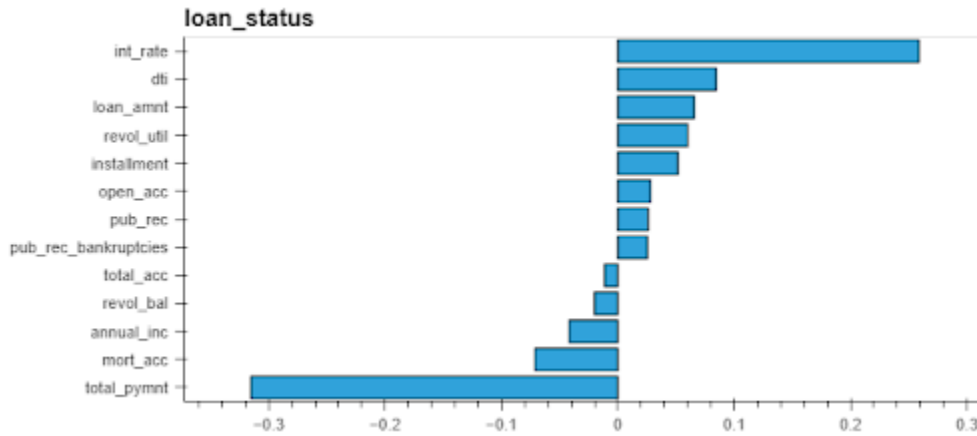


Figure 1: Correlation Coefficient between Target Variable and Various Attributes

3.1 LOGISTIC REGRESSION

The classification report displays precision, recall and f1-score for both the classes 0 and 1 *i.e.*, applicants who will not default and applicants who will default, respectively. Apart from that, the report also shows Macro Average ($0.5 \times \text{score}_0 + 0.5 \times \text{score}_1$) and the weighted average ($w_0 \times \text{score}_0 + w_1 \times \text{score}_1$, where w_0 and w_1 are set according to the number of data points in each class). Accordingly, the Accuracy Score obtained is 97.17% and the Classification Report for the “Train Result” is given in Table 1.

	0.0	1.0	Accuracy	Macro Average	Weighted Average
Precision	0.97	0.99	0.97	0.98	0.97
Recall	1.00	0.87	0.97	0.93	0.97
f1-score	0.98	0.92	0.97	0.95	0.97

Table 1: Train Evaluation Metrics for Logistic Regression

Further, the Accuracy Score is 88.81% and the Classification Report for the “Test Result” is given in Table 2.

	0.0	1.0	Accuracy	Macro Average	Weighted Average
Precision	0.98	0.65	0.89	0.82	0.92
Recall	0.87	0.95	0.89	0.91	0.89
f1-score	0.93	0.77	0.89	0.85	0.90

Table 2: Test Evaluation Metrics for Logistic Regression

The Confusion Matrix for Logistic Regression Model and the ROC Curve for Logistic Regression Model are given in Figures 2 and 3, respectively.

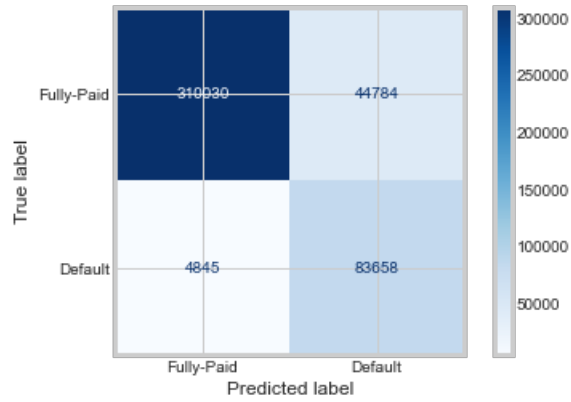


Figure 2: Confusion Matrix for Logistic Regression Model

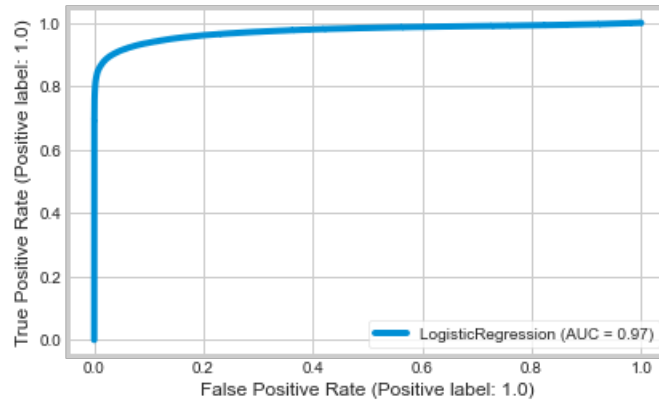


Figure 3: ROC Curve for Logistic Regression Model

3.2 RANDOM FOREST CLASSIFIER

For the Random Forest classifier, the Accuracy Score obtained is 100.00% and the Classification Report for the “Train Result” is given in Table 3.

	0.0	1.0	Accuracy	Macro Average	Weighted Average
Precision	1.00	1.00	1.00	1.00	1.00
Recall	1.00	1.00	1.00	1.00	1.00
f1-score	1.00	1.00	1.00	1.00	1.00

Table 3: Train Evaluation Metrics for Random Forest

Further, the Accuracy Score is 97.11% and the Classification Report for the “Train Result” is given in Table 4.

	0.0	1.0	Accuracy	Macro Average	Weighted Average
Precision	0.97	1.00	0.97	0.98	0.97
Recall	1.00	0.86	0.97	0.93	0.97
f1-score	0.98	0.92	0.97	0.95	0.97

Table 4: Test Evaluation Metrics for Random Forest

The Confusion Matrix for Logistic Regression Model and the ROC Curve for Logistic Regression Model are given in Figures 4 and 5, respectively

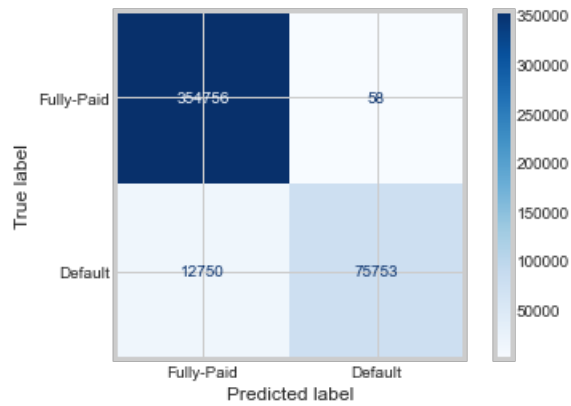


Figure 4: Confusion Matrix for Random Forest Model

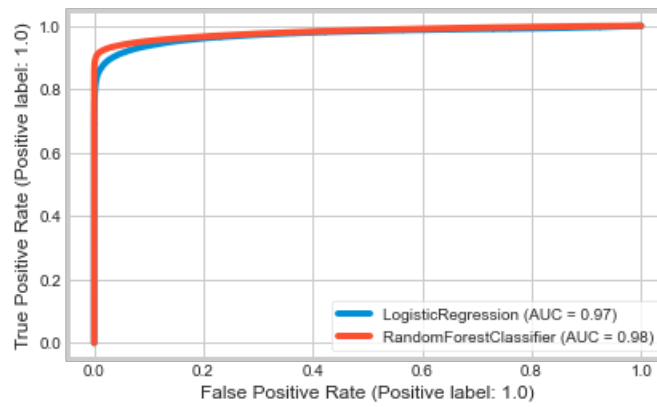


Figure 5: ROC Curve for Random Forest Model

3.3 COMPARISON

Finally, the comparative analysis, based on the AUC are presented in Table 5 and Figure 6. It may be noted from Table 5, that in terms of the AUC score, Logistic Regression performs better than the Random Forest.

	roc_auc_score
Random Forest	0.928
Logistic Regression	0.910

Table 5: AUC Score for Logistic Regression and Random Forest

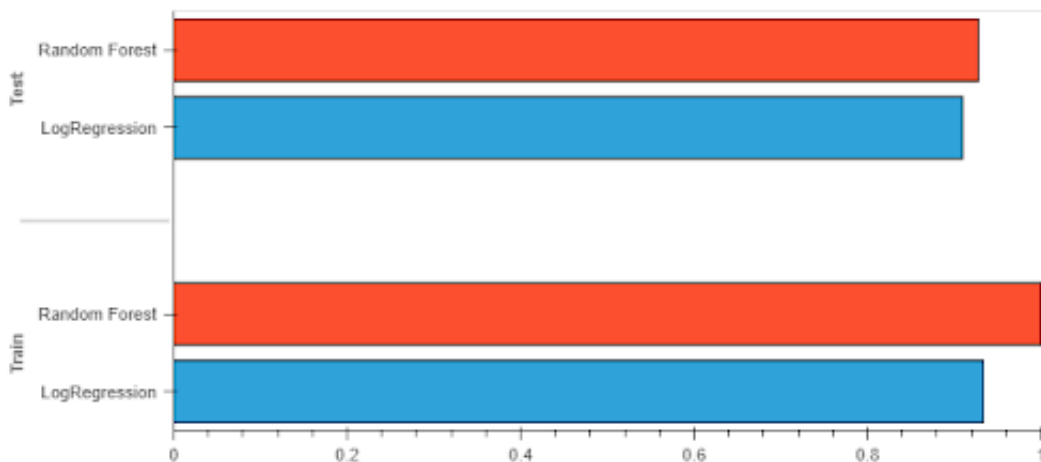


Figure 6: Graph showing Train and Test Accuracy for both the models

3.4 CALCULATING EAD, LGD, EL AND PRICE OF THE CREDIT DERIVATIVE

We begin with the following assumptions in order to calculate EAD and LGD. Recovery rate (which is used to calculate the LGD and the price of the credit derivative) is calculated using the recoveries mentioned for the “Charged Off” clients present in the data. Also, the recovery rate is calculated for each department separately. Pricing of the credit derivative is based upon the theory highlighted in Section 1. Further, the following assumptions are made while pricing the derivative:

- (A) Only 1 payment occurs from protection buyer to protection seller (either at maturity or at the time of default).
- (B) The prevalent (risk-free) interest rates are taken into account while calculating the discounting factor.
- (C) For the calculation of accural amount, we assume that on an average defaults occur right in between of a payment period.
- (D) At the time of the contract writing, the remaining loan amount (including interest rate) is taken as the notional amount and the loan is treated as a bond.

Table 6 shows some (illustrative) debtors with their probability of default, along with EAD, LGD, EL and the price of the credit derivative.

Probability of Default	EAD	LGD	EL	Credit Derivative Price (in bps)
0.09	18330	16727.9	1505.52	0.0469114
0.41	11222.7	10241.8	4199.14	0.156411
0.82	20403.9	18620.6	15268.9	0.447656
0.97	37936.2	34620.6	33581.9	0.38918

Table 6: Table showing EAD, LGD, EL and derivative price for some (illustrative) customers

4 CONCLUSION

The main objective of the paper was the calculation of probability of default for the debtors of a loan lending organization and the consequent design of a credit derivative (*i.e.*, to calculate the premium) thereby hedging against such defaults. Some conclusions based on the results obtained are as follows:

- (A) Based on the performance metrics, the Decision Tree model performs better than the Logistic Regression model.
- (B) Further, when the dependent and independent variables have a high non-linearity and complex relationship, linear models like Logistic Regression will not perform well, since the entire algorithm is based on finding a linear relationship (line of best fit in Logistic Regression) between the variables. In such cases, a tree model (non-linear models) will outperform a classical regression method. Hence, there exists a non-linear relationship between the various attributes and the calculated probability.
- (C) As expected, the calculated derivative premium is higher for loans linked with obligors with high probability of default and short term loans.

REFERENCES

- [1] Wikipedia, Lending Club.
URL <https://en.wikipedia.org/wiki/LendingClub>
- [2] C. Bluhm, L. Overbeck, C. Wagner, Introduction to Credit Risk Modeling, Chapman and Hall/CRC, 2016.
- [3] H. Kim, H. Cho, D. Ryu, Corporate default predictions using machine learning: Literature review, Sustainability 12 (16) (2020) 6325.
- [4] C.-Y. J. Peng, K. L. Lee, G. M. Ingersoll, An introduction to logistic regression analysis and reporting, The Journal of Educational Research 96 (1) (2002) 3–14.
- [5] S. Sperandei, Understanding logistic regression analysis, Biochemia medica 24 (1) (2014) 12–18.
- [6] J. Ali, R. Khan, N. Ahmad, I. Maqsood, Random forests and decision trees, International Journal of Computer Science Issues 9 (5) (2012) 272.
- [7] H. H. Patel, P. Prajapati, Study and analysis of decision tree based classification algorithms, International Journal of Computer Sciences and Engineering 6 (10) (2018) 74–78.
- [8] H. Sharma, S. Kumar, A survey on decision tree algorithms of classification in data mining, International Journal of Science and Research 5 (4) (2016) 2094–2097.
- [9] D. Duffie, K. J. Singleton, Modeling term structures of defaultable bonds, The Review of Financial Studies 12 (4) (1999) 687–720.
- [10] W. Schoutens, J. Cariboni, Levy Processes in Credit Risk, John Wiley & Sons, 2010.
- [11] S. E. Shreve, Stochastic Calculus for Finance II: Continuous-time Models, Springer, 2004.
- [12] J. Garcia, S. Goossens, W. Schoutens, Let's jump together-pricing of credit derivatives: From index swaptions to cpps, Available at SSRN 1358704 (2007).
- [13] T. R. Bielecki, M. Rutkowski, Credit Risk: Modeling, Valuation and Hedging, Springer Science & Business Media, 2013.

- [14] P. J. Schonbucher, *Credit Derivatives Pricing Models: Models, Pricing and Implementation*, John Wiley & Sons, 2003.
- [15] Kaggle, All Lending Club Loan Data.
URL <https://www.kaggle.com/datasets/wordsforthewise/lending-club>
- [16] P. Sonogo, A. Kocsor, S. Pongor, ROC analysis: Applications to the classification of biological sequences and 3D structures, *Briefings in Bioinformatics* 9 (3) (2008) 198–209.