
Personal Loan Financing Model in Brazil - Milestone

Group 2

Salman Aamer, Michael Jennings, Ashwin Ringadoo, Alina Shestiaeva

Abstract

The goal of our project is to build a Machine Learning model capable of predicting whether a given loan should be granted to an individual based on an array of features such as work experience, current earnings, etc. We are planning to look at prior transactional data to determine whether a certain individual taking on a loan will be more likely or not to be placed into a state of default or loan delinquency. We are currently targeting the emerging market of Brazil and focusing on individual consumer loans as we believe this is a segment where there is a strong potential for improvement. The data we have sourced includes a lot of information of individual customer's demographics. It is our objective to not blindly run a regression model to only consider to maximize profit or to minimize the error rate, but to build a model that corrects for any biases that may occur naturally through the regression process. As such we will tackle the biased distribution of loans to create a model where gender, ethnicity and age does not become a determining factor of loan granting and adjusts for any model unfairness.

Introduction and Related Work

There has been significant research on potential methods to evaluate the risk of personal loans (credit scoring). One of the most popular methods in credit scoring is to build the prediction model that estimates the probability of the default of the particular client (Peresetsky, Karminsky, Golovan (2011), Zanin, L. (2020)). However, this kind of model has a lot of factors that needed to be taken into account and treated properly. For instance, if the historical information of a certain client is missing or very rare, the logistic regression shows a very poor predictive performance, and in this case, the log-F prior and ridge regression methods are preferred (Ogundimu, E. O. (2019)). Another important topic is that machine learning algorithms can create discrimination based on protective attributes, such as race, color, religion, gender, and disability. Hardt, Price, Srebro (2016) showed that, for example, Bayes optimal non-discriminating (according to our definition) classifier can intuitively solve this issue. Wattenberg, Viégas, Hardt (2016) indicated that correctly chosen "threshold classifier" (when the bank picks a particular cut-off or threshold, and people whose credit scores are below it are denied the loan, and people above it are granted the loan) can efficiently solve the problem. In our research, we are going to build the probability of the default model considering all discriminatory factors into account.

Data and Methodology

We have decided to make use of the publicly available data from the 13th Pacific-Asia Knowledge Discovery and Data Mining conference. The data includes the credit card applications of Brazilian customers in a one year period. The data will include information on acceptance of application, demographics (age, sex, marital status, etc), income, education. The aim of the project is to evaluate the credit scoring model robustness against performance degradation caused by market gradual changes. Data models for credit card scoring at times are based on profit maximization and can lead

38 to an inherent bias or discrimination of individuals from certain demographics. Our aim is to evaluate
39 how a model can be made without a threshold classifier.

40 0.1 Data Cleaning

41 Prior to the modelling of credit default rates, a substantial data cleaning process is required. The
42 raw dataset has a number of limitations, which include but are not limited to; i.) variables with a
43 large number of missing observations, ii.) variables that have no variation, iii.) variables that are
44 represented numerically but are intended to capture categories, and iv.) variables with invalid values.
45 Table 1 in the Appendix summarises the data cleaning process for each variable.

46 We now briefly describe the procedures in Table 1 in more depth for a number of examples. The
47 variable PAYMENT_DAY indicates which days of the month an individual makes a payment and
48 takes the values 1, 5, 10, 15, 20 or 25. Given that a larger value has no natural interpretation, yet
49 the payment day may have important predictive power, we decide to one-hot encode this feature.
50 RESIDENCE_TYPE takes one of five values, which we also one-hot encode. For AGE, we remove
51 three observations that have values of 3, 7 and 14 as it may be unreasonable for an individual of this
52 age to have a credit card. For the variables that are removed due to missing observations, most of
53 these have over half of the observations missing. AGE is split into six age brackets: 17-25, 26-35,
54 36-45, 46-55, 56-75 and 75 and above. The age brackets are then one-hot encoded. A small
55 number of observations for SEX have entries of 'N', which are removed as we cannot interpret these.

56 0.2 Modelling

57 Our model will aim to look at cross-sectional data on different demographics to ensure that there is
58 equal opportunity. Our constraint will be that of the people that can pay back their loan, the same
59 fraction in each group should be granted a loan. This can be referred to as the true positive rate. Our
60 methods will rely on logistic and linear regression, as well as other machine learning algorithms that
61 may be suitable for the context.

62 We will use the train-test split technique for evaluating the performance of a machine learning
63 algorithm. In order to adjust for any of the potential biases in some of our demographic variables, we
64 will have to first examine the distribution of True Positive rates across several age-groups for example.
65 As such we will know what are the implicit biases included into our model and how we can readjust
66 our training method. It is crucially important as well to determine that a complete accurate model
67 reflecting a 100 percent accuracy is not our objective goal here. The goodness of our model will take
68 into account how the true positive rate will differ across age groups, genders, ethnicity, etc. Splitting
69 the existing data on the training and testing sets, we can evaluate the different metrics that will allow
70 us to check the quality of the model. In addition, we can compare the MSE (or other similar metrics)
71 of our model and the model that were used in the previous similar models from past research.

72 Results

73 The data analysis of our stage involved using multiple Machine Learning Algorithms that would
74 be suitable for a classification problem as presented in this case. The models included a logistic
75 regression, the k-nearest neighbors algorithm (KNN), decision tree algorithm and ensemble learning
76 methods.

77 As our purpose for the model is to ensure that there is no racial or gender discrimination, we started
78 by first checking the level of default in our data sample. Please note that the target variable takes the
79 value of 1, when the individual's debt is considered as bad debt, which would be in the case of a
80 significant delay in making repayments (60 days or more), while it takes the value of 0 in the case
81 when the individual had repaid their debt. The following figures 1,2,3 represent all the results.

82 Using the previous algorithms, two ensemble models were built: a) A model "Mode" that uses a
83 majority voting scheme of the predictions of each algorithm to classify new records. b) A model
84 "Mean" that uses the mean of the probabilities of all the algorithms that the record belongs to the
85 class of interest (i.e., $y=1$). If the mean is higher than the cutoff value (i.e., $p=0.3$), the record is
86 classified in the class of interest. If not, it is classified in the other class. Results are presented in the
87 figure 4.

```

from sklearn.metrics import classification_report
print(classification_report(y_train, logreg.predict(X_train)))

```

	precision	recall	f1-score	support
0	0.74	1.00	0.85	24983
1	0.42	0.00	0.00	8885
accuracy			0.74	33868
macro avg	0.58	0.50	0.43	33868
weighted avg	0.65	0.74	0.63	33868

```

[243] print(classification_report(y_test, logreg.predict(X_test)))

```

	precision	recall	f1-score	support
0	0.74	1.00	0.85	10775
1	0.14	0.00	0.00	3740
accuracy			0.74	14515
macro avg	0.44	0.50	0.43	14515
weighted avg	0.59	0.74	0.63	14515

Figure 1: Classification report of the logistic regression

```

print(classification_report(y_train, knn.predict(X_train)))

```

	precision	recall	f1-score	support
0	0.74	1.00	0.85	24983
1	0.67	0.00	0.00	8885
accuracy			0.74	33868
macro avg	0.70	0.50	0.42	33868
weighted avg	0.72	0.74	0.63	33868

```

print(classification_report(y_test, knn.predict(X_test)))

```

	precision	recall	f1-score	support
0	0.74	1.00	0.85	10775
1	0.00	0.00	0.00	3740
accuracy			0.74	14515
macro avg	0.37	0.50	0.43	14515
weighted avg	0.55	0.74	0.63	14515

Figure 2: Classification report of the KNN

Our initial analysis was performed on the set of data after removing discriminatory variables including gender, age, marital status. The variables included were standardized income (proxy for being employed), other incomes (proxy for having sufficient funds for repayment in the case of being unemployed), quantity of cars (proxy for personal asset value), months in residence (proxy for being a local resident), possession of credit card (proxy for having a credit history), self-reported employer name (proxy for being employed and ability to perform background checks), and type of product applied (encoding is unknown, however we believe it will help our model in choosing the correct parameters).

The analysis performed with the various models led us to select the logistic regression as our preferred model. We have used Grid Search to help choose the optimum parameters. The classification accuracy on the test data for the models used are presented below. As the figures below indicate, the logistic regression leads us to having the greatest precision on both of the target variable values (0 and 1). Our analysis is in line with the literature presented above which highlights the use of a logistic regression in the case of sufficient historical client information (Ogundimu, E. O. (2019)).

We will be analyzing the data using the possible discriminatory variables (age, gender, marital status) to evaluate the true positive rates from our model on these individual variables to evaluate if there is a bias and will then correct for this bias.

```

▶ print(classification_report(y_train,tree.predict(X_train))
⊖

```

	precision	recall	f1-score	support
0	0.74	1.00	0.85	24983
1	0.00	0.00	0.00	8885
accuracy			0.74	33868
macro avg	0.37	0.50	0.42	33868
weighted avg	0.54	0.74	0.63	33868

```

/usr/local/lib/python3.7/dist-packages/sklearn/metrics/_cl
_warn_prf(average, modifier, msg_start, len(result))

[ ] print(classification_report(y_test,tree.predict(X_test)))

```

	precision	recall	f1-score	support
0	0.74	1.00	0.85	10775
1	0.00	0.00	0.00	3740
accuracy			0.74	14515
macro avg	0.37	0.50	0.43	14515
weighted avg	0.55	0.74	0.63	14515

```

/usr/local/lib/python3.7/dist-packages/sklearn/metrics/_cl
_warn_prf(average, modifier, msg_start, len(result))

```

Figure 3: Classification report of the decision tree

```

[391] classificationSummary(total_valid.actual, total_valid.Mean)

```

Confusion Matrix (Accuracy 0.7007)

	Prediction	
Actual	0	1
0	9790	985
1	3360	380

```

▶ classificationSummary(total_valid.actual, total_valid.Mode)

```

Confusion Matrix (Accuracy 0.7423)

	Prediction	
Actual	0	1
0	10775	0
1	3740	0

Figure 4: Classification report of the ensemble model

References

- 106 1) Peresetsky, Anatoly A., Alexandr A. Karminsky, and Sergei V. Golovan (2011). "Probability of
107 default models of Russian banks." *Economic Change and Restructuring* 44.4: 297-334.
- 108 2) Ogundimu, E. O. (2019). Prediction of default probability by using statistical models for rare
109 events. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 182(4), 1143-1162.
- 110 3)Hardt, M., Price, E., Srebro, N. (2016). Equality of opportunity in supervised learning. №
111 1610-02413.
- 112 4)Wattenberg, M., Viégas, F., Hardt, M. (2016). Attacking discrimination with smarter machine
113 learning. *Google Research*, 17.
- 114 5)Zanin, L. (2020). Combining multiple probability predictions in the presence of class imbalance to
115 discriminate between potential bad and good borrowers in the peer-to-peer lending market. *Journal*
116 *of Behavioral and Experimental Finance*, 25, 100272.
- 117 **Github Repository:** <https://github.com/salmanaamer/6.862.Spring.2021.Group.2.git>

Table 1: Summary of Data Cleaning

Data Limitation	Variables	Solution
Categorical variables with more than two categories	RESIDENCE_TYPE PAYMENT_DAY PRODUCT	One-hot encoding
Variables with invalid values	SEX	Remove observations
Unreasonable values	AGE MARITAL_STATUS	Remove observations
Constant values	EDUCATION_LEVEL CLERK_TYPE QUANT_ADDITIONAL_CARDS FLAG_MOBILE_PHONE	Drop Variable
Too many missing observations	MATE_EDUCATION_LEVEL APPLICATION_- SUBMISSION_TYPE RESIDENCIAL_PHONE- _AREA_CODE PERSONAL_ASSETS_VALUE PROFESSIONAL_STATE PROFESSIONAL_CITY PROFESSIONAL_BOROUGH PROFESSIONAL_PHONE- _AREA_CODE MATE_PROFESSION_CODE FLAG_HOME_ADDRESS- _DOCUMENT FLAG_RG FLAG_CPF FLAG_INCOME_PROOF FLAG_ACSP_RECORD	Drop variable
No clear interpretation of encoding	NACIONALITY	Drop variable
Duplicate variable	QUANT_SPECIAL- _BANKING_ACCOUNTS	Drop variable
Continuous variables used for discrimination testing	AGE	Divide into 6 age buckets and one-hot encoded
Geographic variables	STATE_OF_BIRTH CITY_OF_BIRTH RESIDENCIAL_STATE RESIDENCIAL_CITY RESIDENCIAL_BOROUGH PROFESSION_CODE PROFESSIONAL_ZIP_3 RESIDENCIAL_ZIP_3	We will keep one of these variables but are in the process of encoding it
Similar variables	PERSONAL_MONTHLY- _INCOME OTHER_INCOME	Summed and then standardised
In progress	OCCUPATION_TYPE	Encoding in progress