

PROBLEM LOAN PREDICTION

DEFINITION

Project overview

Banks, and other financial institutions, make their profits by lending money (www.investopedia.com). Every time a loan is issued the institution takes a risk. Risks include the fact that some borrowers, or customers, will not return the money as agreed during the lending process. Lack, or delay, payments results in “Problem loans,” or loans where the banks lose money as they cannot recover the original amount, or have to invest additional money during the repaying process (www.iedunote.com). According to the US Department of Treasury, banks deal with problem loans by: renewing or extending the loan terms, extending /adding credit, restructuring the loan, or foreclosure. Regardless of the direction taken, during this process, known as loan workout, the bank is forced to choose the alternative where the recovery is maximized and the risk and expenses are minimized (www.occ.treas.gov). In these circumstances, machine learning techniques and historical data can provide additional information during the initiation of the loan process in order to reduce the risk by identifying customers likely to generate a problem loan.

In this project, an artificial neural network classifier is built to predict problem loans on data from the Lending Club dataset (<https://www.lendingclub.com/info/download-data.action>) corresponding to quarters 1 and 2 of 2017. The model is built in Python 3 using Keras with Tensorflow.

Problem Statement

The purpose of this project is to generate a solution to predict whether or not the loan application of a Lending Club customer will be fully paid loan or generate a problem loan for the lender. The proposed solution will be a 3-layer artificial neural network that will be trained and tested on the available information for 2017.

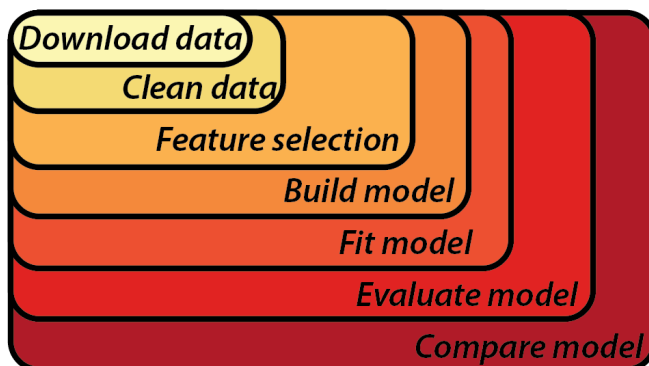


Figure 1 Schematics organization of the procedures taken in the creation of a predictive model for problem loans

To build the solution, the data will be obtained from the Lending Club statistics website (<https://www.lendingclub.com/info/download-data.action>), loaded into Python, cleaned, and submit into the model (Figure 1). The model will be tested against subsets of the original dataset and the results compared to two different models linked to a Kaggle dataset that compiles data from 2007 to 2015 (<https://www.kaggle.com/wendykan/lending-club-loan-data>).

Metrics

Originally, the proposal for this project included the use of the F-score to evaluate the model and compare to other models on similar datasets. The F-score provides a balance between precision and recall, and therefore complements these two measures, and reduce the accuracy paradox. As comparisons with other models were executed (See below) both kaggle models failed to identify problem loans and therefore the F-score was not calculated. To overcome this situation, accuracy (total of correct predictions divided total number of loans) was used to determine possible differences between the proposed model and the other models.

ANALYSIS

The proposed solution here is divided into seven phases as depicted in figure 1.

Download the data

The data comes from the Lending Club statistics, it includes quarters 1 and 2 of 2017 (<https://www.lendingclub.com/info/download-data.action>). The data comes in two files, one for each quarter that are loaded into Python separately and merge into a single data frame. The final data contains a total of 139,919 loans with 145 columns of information. The information in the columns include information such as: loan type, amount, dates, and personal data on the credit history of the customer as well as the stats of the loan. The specific descriptions on each of the 145 columns is in: LCDataDictionary.xlsx.

Data exploration, cleaning and visualization

An initial exploration of the data resulted in two obvious problems that need to be addressed:

- a. There are text descriptions of the data in the last rows of the id column (id column is empty since it is personal information). These cells were removed after sorting the data and deleting the last 8 rows. This procedure resulted in 139,311 loans.
- b. There are many cells with NaNs. Some columns have no information at all, but are filled with NaNs. These columns were identified and removed after counting NaNs. Initially reducing the number of columns to 142 and then to 36, maintaining the number of NaN per column to a maximum of 10.

Feature selection

Target feature

The target feature corresponds to `loan_status`. There are nine different classes of loans, these can be grouped into three main categories:

- *Current loans* do not show problems, payments are current and for this reason were removed from the dataset.
- *Fully paid* loans are those that have been paid in full. It is not clear whether they were paid on time or not, but the assumption is that they did not show problems or delays before payment.
- *Problem loans* are all loans that have a problem, including Default, Grace Period, Do not meet credit policy, and/or Charged off.

Fully paid and problem compromise 58,006 loans.

Predicting features

The other 35 features (excluding `status_loan`) were classified into: non-numeric/categorical and numeric features.

There are 18 non-numerical variables. From these, `zip_code`, `sub_grade`, and `initial_list_status`, were removed. Data in `zipcode` included only 877 of the 43,000 in the USA and the information they may provide to the model should be contained in the `state` column. Similarly, `subgrade` is a subcategory of `grade` and `initial_list_status` is not clearly described.

Interest rate (`int_rate`) and employment length (`emp_length`) had string characters in most of the cells. These characters were replaced or removed where necessary and the columns converted into a numeric variable.

Finally, two columns, issued date (`issue_d`) and last credit check (`last_credit_pull_d`) have date information. Both of these columns are formatted as “Month-year”. This data was converted into delta time by calculating the time in days from to the last date in the data set. After this transformation, both variables became numeric.

After processing the abovementioned columns, the dataset contains 11 categorical features All features are more or less well represented across the two categories (Fully paid and Problem loans) of the loan status variable (Figure 2).

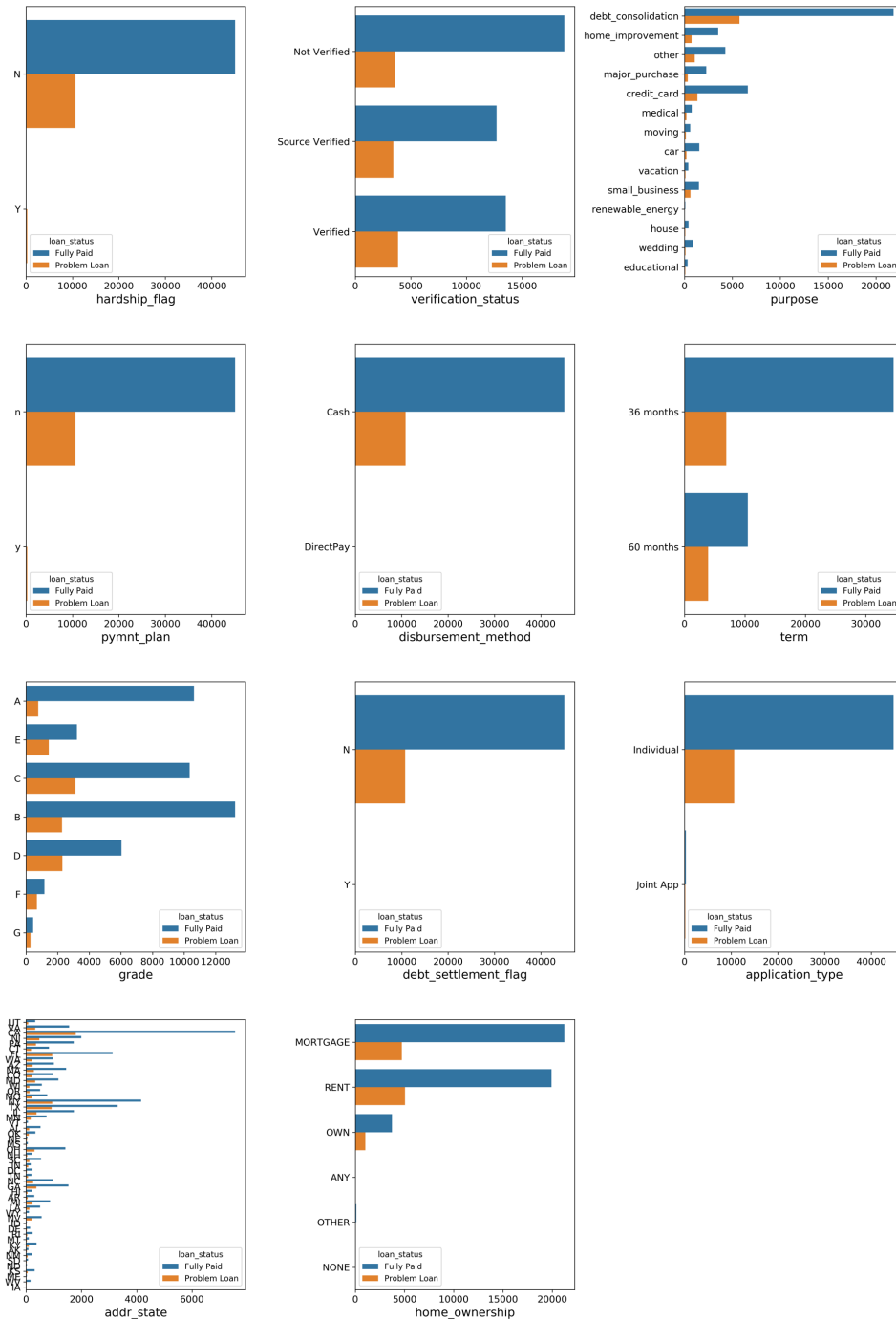


Figure 2. Bar plots showing the counts of fully paid and problem loans for each category in the categorical variables from the 2017 first two quarters of data.

The original number of numeric variables was 17, but with the conversion of categorical variables, this number increased to 21.

Of the numeric columns, policy_code had a single value (1). Installment, funded_amnt, funded_amnt_inv, and out_prncp_inv corresponded to similar types of data (loan_amnt and out_prncp – Table 1). These 4 features were removed.

Table 1 Number of unique observations for each of the numeric variables in the 2017 loan dataset.

<i>Numeric feature</i>	<i>No. of unique values in the feature</i>
loan_amnt	1208
revol_bal	26976
out_prncp	3447
out_prncp_inv	3460
total_pymnt	55342
total_pymnt_inv	52965
total_rec_late_fee	2627
recoveries	4976
collection_recovery_fee	2889
last_pymnt_amnt	47818
policy_code ***	1
total_rec_int	48546
total_rec_prncp	11752
funded_amnt	1266
int_rate	420
installment	19462
emp_length	11
issue_d	6
funded_amnt_inv	9307
annual_inc	6415
last_credit_pull_d	124

A correlation analysis between numerical values did not indicate any concerning correlation between any pair of variables (Figure 3). Only total_rec_prncp showed higher correlations with total_payment ($r^2 = 0.97$) and total_paymnt_inv ($r^2 = 0.95$ Figure 3). All other correlations showed an r^2 lower than 0.80.

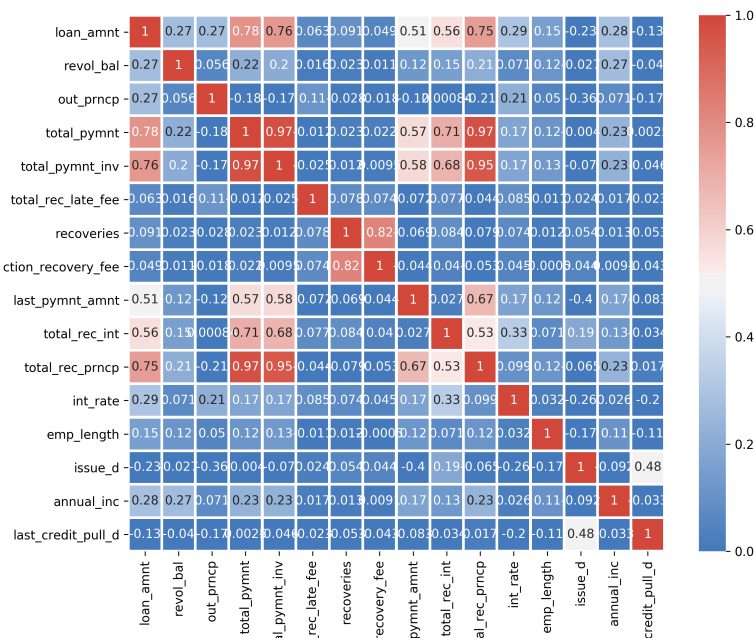


Figure 3 Correlations between numerical variables in the 2017 Lending Club dataset.

A simple boxplot visualization of the numeric variables showed the presence of outliers in most of them (Figure 4 - top).

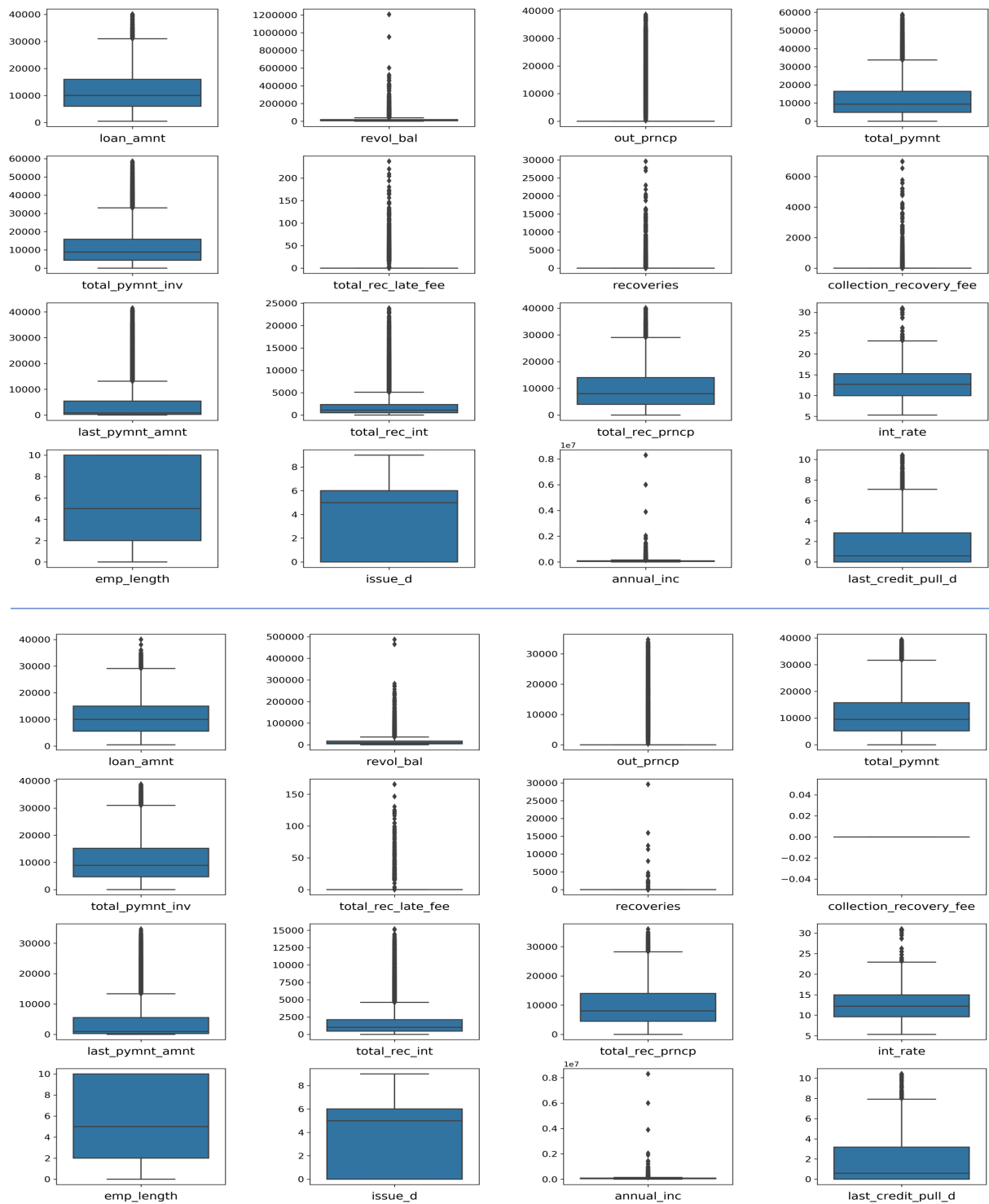


Figure 4 Boxplots of the numeric variables before (top) and after (bottom) outlier removal.

Outlier analysis using the 1st and 3rd percentile revealed 8383 rows with data identified as outliers across all columns. These data points were eliminated. Loans with outliers on fewer columns were kept in the dataset. After this step, the number of loans in the dataset is 47,505 with 28 features.

Finally, the 11 categorical variables were encoded using LabelEncoder from scikit-learn, loan_status was also encoded. This is necessary for the neural network to work, as it does not take categorical variables.

Algorithms and Techniques

The proposed classifier model, an artificial neural network, consists of:

- one input layer
- one hidden layer
- one output layer.

All layers are regular densely-connected layers uniformly initiated. The input and hidden layers use the relu activation function with 14 units. The output layer is activated by a sigmoid function with a single unit.

The network compiler uses the adam function as optimizer, binary_crossentropy for the loss, and set the accuracy as the metrics for the model.

Benchmark model

There are a few methods that have used data sets from the Lending Club statistics site. Of these two will be used as benchmark models:

- **Kaggle-1** (<https://www.kaggle.com/mina20/building-a-neural-net-to-predict-default/notebook>)
- **Kaggle-2** (<https://www.kaggle.com/gagrawal/neural-net-with-keras>)

This used neural network for the predictions, kaggle-1 uses tensorflow directly and kaggle-2 uses keras and tensorflow. Neither report metrics, but after running them here (See below), they have an accuracy of 94.47 and 93.0% respectively. In addition to running the model and estimating the accuracy, the F-score will also be calculated with the predictions of these models.

METHODOLOGY

Data preprocessing

In order to prepare the data for the model, it was split into the target variable (loan_status) and the predictors (features). Then these portions are divided into training (75%), validation (20%), and test (5%) subsets.

The features were standardized using StandardScaler which removes the mean and scales the data to unit of variance. This step is necessary to prevent bad behaviors in the model particularly from features that are not normally distributed.

Implementation

The model is built in keras using TensorFlow backend, after compilation, it will be trained with the training subset. The resulting fitting will provide the accuracy of the model. The model will then be then validated in two ways:

- Cross-validation using the training set
- Validation using the validation set

Once the model has been validated, it will be tested on the test data. The predictions from the test run will be used to determine accuracy and F-score of the model. To compare the current model with the benchmark models, the 2007-2015 data used by the kaggle models will be used.

RESULTS

Model evaluation and validation

After compiling the model, it was fitted with the training data set, run for 10 epochs with a batch size of 10. This initial run had an accuracy of 99.74% and a loss of 0.0093. The cross validation of the model using a k-fold of 10 returned an average accuracy of 99.65% +/- 0.15 and the validation with the validation subset was 99.6913% accurate with only 0.0118 loss.

These high accuracies might be the result of overfitting, even though they were already tested on different subsets. However, as a precaution the classifier was tested on data from the 4th quarter of 2016 (see notebook 2016_pred.ipynb). This test resulted was of 98.49% (20,005 loans) accurate. The model was also tested in a dataset that included loans from 2007-2015 and is compiled in kaggle (<https://www.kaggle.com/wendykan/lending-club-loan-data>). In this larger dataset, the model showed an accuracy of 99.82% (see notebook 2007-2015_pred.ipynb). It is very unlikely that the model is over fitted given that all datasets are different.

Model optimization

An optimization of the model using grid search that iterates through different parameters suggested to change the model by adding 10 extra epochs.

The model with 20 epochs resulted in the same accuracy of 99.74% as the initial run. The validation was also 99.74%.

Implementation

The optimized classifier on the test set had an accuracy of 99.72% and an F-score of 0.9871. It only misclassified 43 loans out of the 11,877 loans (Figure 5).

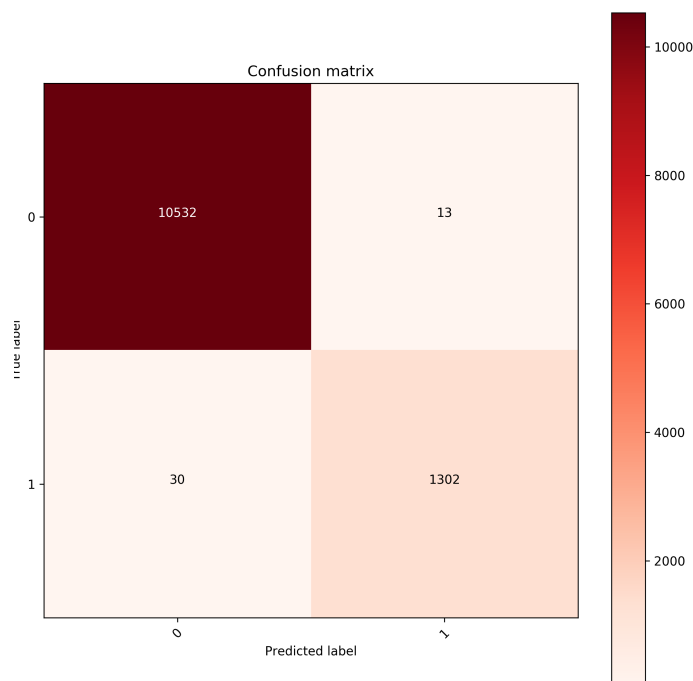


Figure 5. Confusion matrix of the 2017 dataset predictions.

Comparison with benchmark models

The proposed ANN classifier was fitted with the data from 2007-2015 as to not to change the kaggle models to use the 2017 dataset. Kaggle-1 and kaggle-2 predict default loans on Lending Club data from 2007 to 2015 and both also use neural networks as their preferred method.

Kaggle-1 uses tensorflow with a DNNRegressor on 78 features and 37,379 loans. The original script only reports the loss during fitting with a value of 0.0253. The measured accuracy (number of correct predictions over total number of predictions) was 94.47%. This model, however, was not able to identify *problem loans*; therefore, it was not possible to determine the F-Score (Table 2, Figure 6A).

Kaggle-2 uses keras on tensorflow with a 3-layer classifier network and 68 features on 1,000 loans. The loss during fitting was 1.4029 and the testing accuracy 91.2%. This model also failed to predict *problem loans*, and the F-Score was not determined (Table 2, Figure 6B).

Table 2 Summary table with the comparisons between the proposed model and the models from two models uploaded in kaggle.

	<i>Kaggle-1</i>	<i>Kaggle-2</i>	<i>ANN Proposed here</i>
<i>Model</i>	DNN Regressor	Classifier 3 layers	Classifier 3 layers
<i>Library</i>	Tensorflow	Keras-tensorflow	Keras-tensorflow
<i>No. loans tested</i>	37,379	1,000	97,060
<i>No. features</i>	78	68	23
<i>Loss during fitting</i>	0.025	1.12	0.01
<i>Accuracy on test dataset</i>	94.47	93.0	97.30
<i>F-Score</i>	Did not predict Problem loans 1 in Figure 6A	Did not predict Problem loans 0 in Figure 6B	0.9058

The ANN model proposed here is also based on keras and tensorflow with 3-layers, but it uses 23 features and 97,060 loans for testing. This model has an accuracy of 99.78%, loss of 0.0119 during fitting, and predicted both problem and paid loans with an F-Score of 0.9058 with an accuracy of 97.30% (Table 1, Figure 6 C).

CONCLUSION

Free-visualization

Figure 6 shows three different confusion matrixes showing the accuracy of the model as compared to the kaggle models.

Reflection

The proposed model outperforms other methods by 5-7 points in terms of accuracy and predicts both categories. The outperformance could be the result of:

- Cleaning data
- Labeling of the categorical features.
- The number of selected features.
- The standardization of the features before the model.

Kaggle-1 uses a different structure to that in the current model, making the comparison slightly difficult. On the other hand, kaggle-2 and the ANN proposed here have the same 3-layer

structure but the labelling is different, and in kaggle-2 the number of selected features is higher, and it lack feature standardization. Combined these differences may account for the differences in accuracy seen in the models. Models using traditional machine learning algorithms result in lower accuracies, as shown by at least one project that used R and several models to predict problem loans in a dataset from the Lending Club (https://rstudio-pubs-static.s3.amazonaws.com/203258_d20c1a34bc094151a0a1e4f4180c5f6f.html).

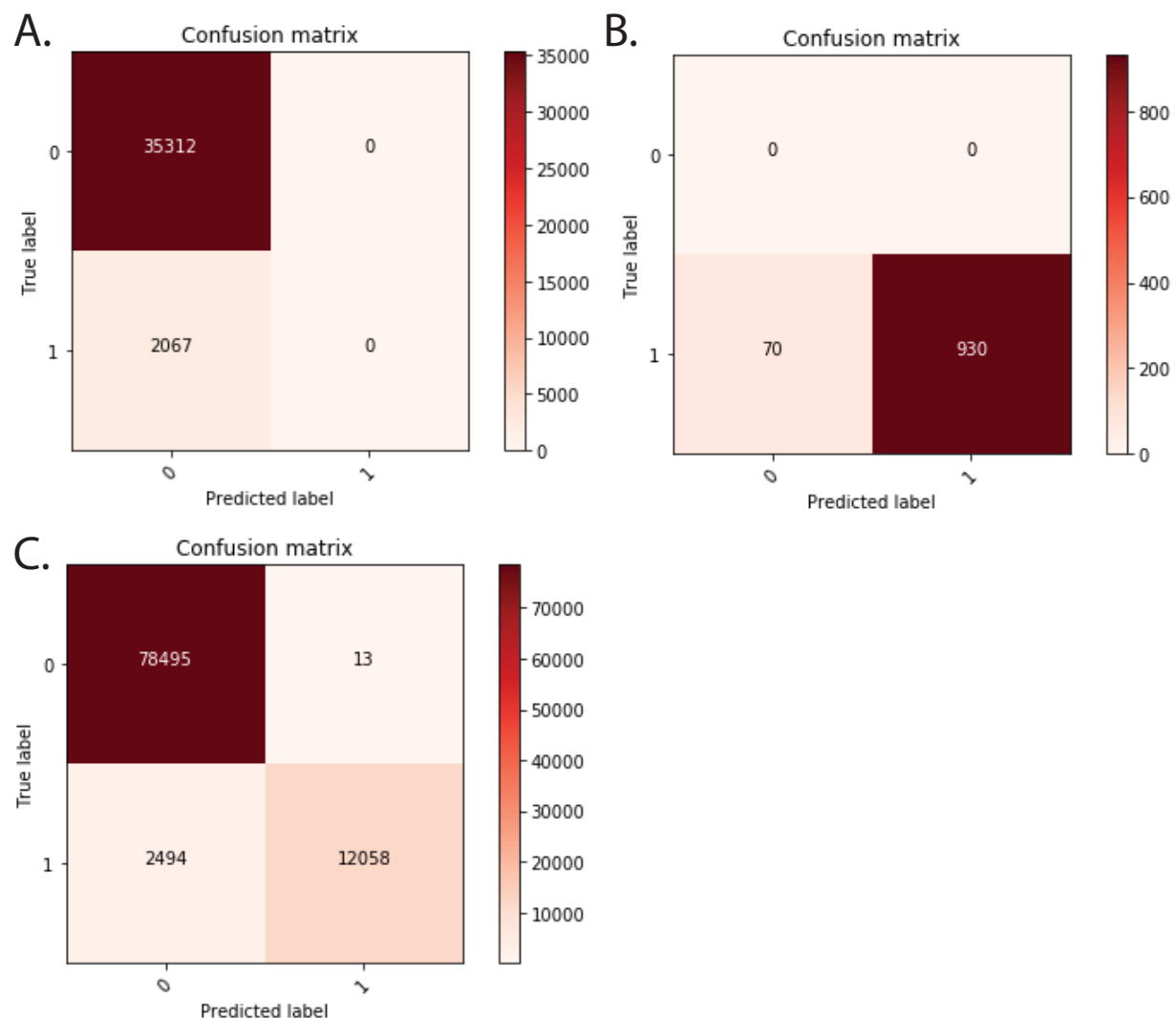


Figure 6 Confusion matrices for kaggle-1 (A), kaggle-2 (B), and the ANN proposed in this project (C) using the 2007-2015 dataset.

Reflection

The model proposed here is highly effective, with accuracies higher than 99.5% at predicting both paid and problem loans across different datasets from the Landing Club. It outperformed

other models, likely as a results of improved data cleaning, encoding, and standardization, as well as the use of a smaller set of features.

Improvement

With accuracies of 99.5% and larger and F-Scores bordering 0.98, it is difficult to come up with an improvement for the model. Rather it will be interesting to:

- Add more features and determine how the accuracy is affected.
- Test this model in datasets other than those from the Lending Club.