

Financial DATA SCIENCE

Group Project 5

Carlos Caballos - 20220528

Frederico Bravo - 20221231

Neftali Herculano - 20200732

Nelson Lima - 20221539

Riccardo Gurzu - 2022123

<i>CONTENTS</i>	2
-----------------	---

Contents

1 Introduction	3
2 Primary Statistical Analysis	4
3 Data Pre-processing	5
3.1 Missing Values	6
3.2 Outliers	7
3.3 Feature Engineering	7
4 Modelling	9
4.1 Feature Selection	9
4.2 Logistic Regression	10
4.3 Random Forest	12
4.4 Deep Neural Network	13
5 Conclusion	15
References	16
Code	17

List of Figures

1 Target Feature Distribution	5
2 Variables boxplot	7
3 Correlation Heatmap	8
4 Correlation with the Target Feature 'risk'	10
5 Logistic Regression ROC Curve	12
6 Random Forest ROC Curve	13
7 Deep Neural Network ROC Curve	14

List of Tables

1 Descriptive statistics for the Numerical Features	5
2 Threshold for the Outliers	8
3 Confusion Matrix Logistic Regression	11
4 Performance measures Logistic Regression	11
5 Confusion Matrix Random Forest	13
6 Performance measures Random Forest	13
7 Confusion Matrix Deep Neural Network	14
8 Performance measures Deep Neural Network	14

Abstract

In this study, we were challenged to develop a data modelling for a credit dataset to understand the customer's loan profiles better. To do so, we went through a primary statistical analysis of the credit dataset that considered outliers' detection and treatment of missing values to understand better the data given.

Secondly, we also developed three ML models (Logistic Regression, Random Forest and Deep Neural Network) to predict the probability of default and evaluated the model's results. We also tested alternative models and compared our results. We decided to use the R version 4.2.0 (2022-04-22) for calculations and \LaTeX to prepare the final report.

The findings of this study show that the Random Forest methodology represents a more reasonable coefficient of determination and a lower false-negative rate than other models.

1 Introduction

For the financial sector, particularly banks, credit risks represent a significant risk. It occurs when a firm or the debtor does not fulfil debt obligations individually. That is to say, that the probability of a lender not having received principal and interest on debt can also be taken into account. This is why board members and senior management must establish and demonstrate a clear understanding of the measurement of credit risk in their institutions. Although accounting firms and bank regulators have been going in a similar direction regarding risk modelling, they each develop their own models by relying on the best instruments at their disposal to seek reasonable values and attempt to predict potential losses. Today, banks and financial institutions are using Big Data and Machine Learning Models to help predict credit risk, which is the process of predicting who to lend money to and how much to extend credit, and credit monitoring, which tracks the borrower's credit behaviour over time. Financial institutions are also using ML to detect fraud and prevent loss.

The history of credit risk assessment and its intersection with machine learning reflects an evolution driven by the increasing complexity of financial landscapes and the growing availability of diverse data sources. Traditionally, credit risk evaluation relied on rule-based systems and statistical models with limited capacity to adapt to dynamic market conditions. Over the past few decades, the surge in computing power and the proliferation of digital data have catalyzed a paradigm shift towards machine learning in credit risk management. The use of neural networks, decision trees, and ensemble methods has become more prevalent, allowing for a more nuanced understanding of risk factors and patterns. Notably, the 2008 financial crisis underscored the shortcomings of conventional risk models, prompting a renewed focus on incorporating advanced machine learning algorithms capable of capturing nonlinear relationships and identifying

hidden risks. Since then, researchers and financial institutions have continued to refine and expand machine learning applications in credit risk assessment, aiming to enhance predictive accuracy and anticipate default events with greater precision in an ever-changing financial landscape.

In conclusion, the convergence of credit risk assessment and machine learning represents a pivotal juncture in the history of financial decision-making. As institutions seek to navigate the intricate web of factors influencing creditworthiness, the adoption of machine learning methodologies stands as a testament to the ongoing quest for more robust and adaptive risk management tools. This symbiotic relationship between data-driven algorithms and the historical context of credit risk not only addresses the limitations of traditional models but also reflects a commitment to staying ahead of emerging challenges in the global financial ecosystem. By harnessing the power of machine learning, stakeholders in the financial industry aim not only to accurately predict defaults but also to foster a more resilient and responsive approach to credit risk assessment, ultimately shaping a landscape where precision, agility, and informed decision-making are paramount. This is exactly what we will be doing during this project using the a training and test dataset to develop a model that will be able to predict if the credit would or not default and understand the risk of it.

2 Primary Statistical Analysis

An Exploratory Data Analysis is the first step in any data analysis project. In the present project, 310704 observations on the training dataset and 20600 on the test dataset need to be explored in three different ways:

- i. Summarizing the dataset using descriptive statistics;
- ii. Visualizing the dataset using charts;
- iii. Identifying missing values.

Before starting to perform statistical models, the three actions will enable us to get a sense of how values are distributed and whether there are any problems.

In the dataset, we have two types of variables, categorical and numerical:

- Categorical/factor variables: *term*, *grade*, *emp_title*, *emp_length*, *home_ownership*, *verification_status*, *issue_d*, *purpose*, *addr_state*, *earliest_crline* and *loan_status*.
- For the numerical values, we present the descriptive stats below:

On Figure 1 it is possible to understand that we have more observations, on the train dataset, classified with 0 on risk than with 1, meaning that majority of the loan didn't default. To be more precise, 207036 observations didn't default while 103668 did default.

Variable	n	mean	sd	median	min	max	range	skew	kurtosis	se
loan_amnt	310704	15518.61	9196.53	14000.00	1000.00	40000.00	39000.00	0.70	-0.27	16.50
funded_amnt	310704	15518.61	9196.53	14000.00	1000.00	40000.00	39000.00	0.70	-0.27	16.50
funded_amnt_inv	310704	15511.82	9195.22	14000.00	725.00	40000.00	39275.00	0.70	-0.27	16.50
int_rate	310704	12.57	4.70	11.99	5.32	30.99	25.67	0.75	0.27	0.01
installment	310704	452.84	264.51	387.55	14.77	1618.24	1603.47	0.95	0.55	0.47
annual_inc	310704	80539.98	92862.94	67200.00	0.00	9757200.00	9757200.00	46.99	3628.71	166.60
dti	310556	19.02	12.34	18.19	-1.00	999.00	1000.00	24.10	1587.69	0.02
delinq_2yrs	310704	0.34	0.92	0.00	0.00	21.00	21.00	5.39	49.12	0.00
inq_last_6mths	310703	0.61	0.89	0.00	0.00	5.00	5.00	1.70	3.12	0.00
open_acc	310704	11.88	5.79	11.00	0.00	81.00	81.00	1.33	3.40	0.01
pub_rec	310704	0.25	0.67	0.00	0.00	86.00	86.00	17.17	1369.85	0.00
revol_bal	310704	16052.73	23228.24	10590.00	0.00	1044210.00	1044210.00	9.62	192.60	41.67
revol_util	310491	48.50	24.81	48.00	0.00	182.80	182.80	0.06	-0.80	0.04
total_acc	310704	24.92	12.31	23.00	2.00	176.00	174.00	1.05	2.15	0.02
out_prncp	310704	723.86	3602.73	0.00	0.00	40000.00	40000.00	6.17	42.23	6.46
total_pymnt	310704	13430.48	10010.94	10741.86	0.00	59808.26	59808.26	1.11	0.73	17.96
risk	310704	0.33	0.47	0.00	0.00	1.00	1.00	0.71	-1.50	0.00

Table 1: Descriptive statistics for the Numerical Features

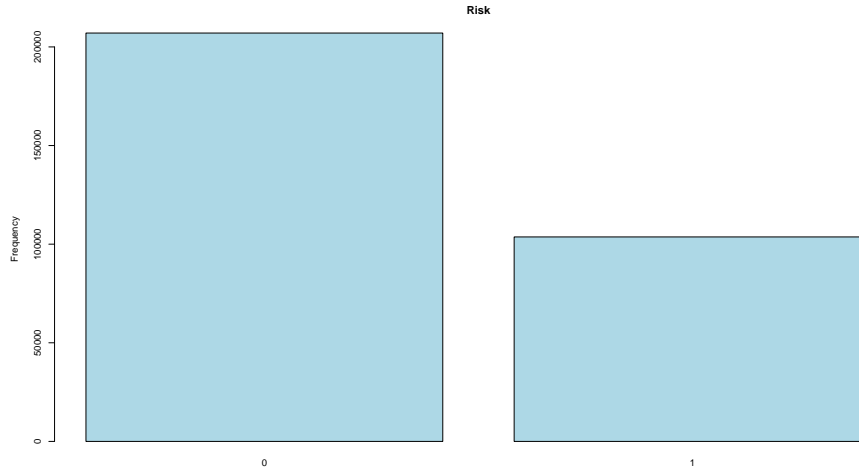


Figure 1: Target Feature Distribution

3 Data Pre-processing

Outliers can lead to vague or misleading predictions while using machine learning models. They decrease the mathematical power of statistical models, and thus, the output of the models becomes unreliable. Removing and modifying outliers using statistical detection techniques, such as Z-score, density-based spatial clustering, regression analysis, proximity-based clustering, and IQR scores, is a widely followed method. We can also use visual detection, such as box plots, to identify outliers.

Missing values are empty cells, rows, and columns that we often see in a dataset. They make the dataset inconsistent and unable to work on. Many

machine learning algorithms return an error if parsed with a dataset containing null values. Detecting and treating missing values is essential while analysing and formulating data for any purpose.

3.1 Missing Values

Treating missing values and outliers is an essential step of data cleansing and preparation and should be one of the first operations done on a dataset. Only after these steps are administered is the data considered usable to build models and take insights from. We found the following missing values in the training dataset while performing remedies (we need to write down the methods we used):

- *emp_title*: need correct for missing values of 29563 observations: we considered as the most conservative approach (“Unemployed”);
- *emp_lenght*: need correct for missing values of 22615 observations: we considered as the most conservative approach (< 1 year);
- *dti*: need correct for missing values of 148 observations: we decided to apply the mean;
- *ing_last_6mths*: need correct for missing values of 1 observation: we decided to apply the mean;
- *revol_util*: need correct for missing values of 213 observations: we decided to apply the mean.

And we found the following missing values in the validation (test) dataset:

- Dropping the row where *loan_amnt* has missing values since 180 rows are not correctly separated;
- *emp_title*: need correct for missing values of 1224 observations: we considered as the most conservative approach (“Unemployed”);
- *emp_lenght*: need correct for missing values of 1210 observations: we considered as the most conservative approach (< 1 year) ;
- *dti*: need correct for missing values of 18 observations: we decided to apply the mean;
- *revol_util*: need correct for missing values of 17 observations: we decided to apply the mean.

3.2 Outliers

Moving to the outliers, the following boxplot shows a general view of the dispersion of the variables used in the training dataset:

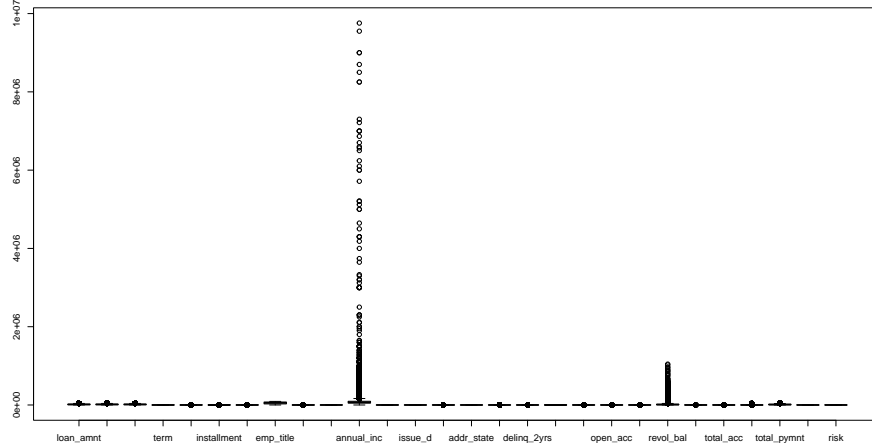


Figure 2: Variables boxplot

We decided to exclude from this checks by outliers the variables *emp_title*, *issue_d* and *earliest_cr_line* due to the high number of levels, which would make it complicated and even impossible for some of them to check for the outliers.

To deal with the outliers, we defined some thresholds for the numerical variables based on the boxplots and histograms of each feature, and all of the observations that didn't comply with that threshold will be automatically excluded. The thresholds are given in Table 2.

We applied the threshold to both the train and test dataset. By this step we conclude the pre-processing part of the project.

3.3 Feature Engineering

The feature engineering process included correlation analysis, data type conversions, mapping of categorical variables to numerical values, and the creation of new features, these steps were crucial in order to improve the model performance and providing meaningful insights from the data.

The initial step was the Numerical Columns Identification, where numerical columns were identified from the *train_dataset* using the *supply* function with *is.numeric*, which helped in focusing on quantitative data for correlation analysis. Then the second step was the creation of a correlation matrix for the numerical columns and then visualized by using a heatmap (Figure 3), this helped to understand the relationships between different numerical features.

Variable	Threshold
int_rate	25
annual_inc	300000
dti	50
delinq_2yrs	5
inq_last_6mths	3
open_acc	30
pub_rec	3
revol_bal	60000
revol_util	110
total_acc	70
out_prncp	10000
total_pymnt	50000
credit_line_duration	20000
income_to_loan_ratio	40

Table 2: Threshold for the Outliers

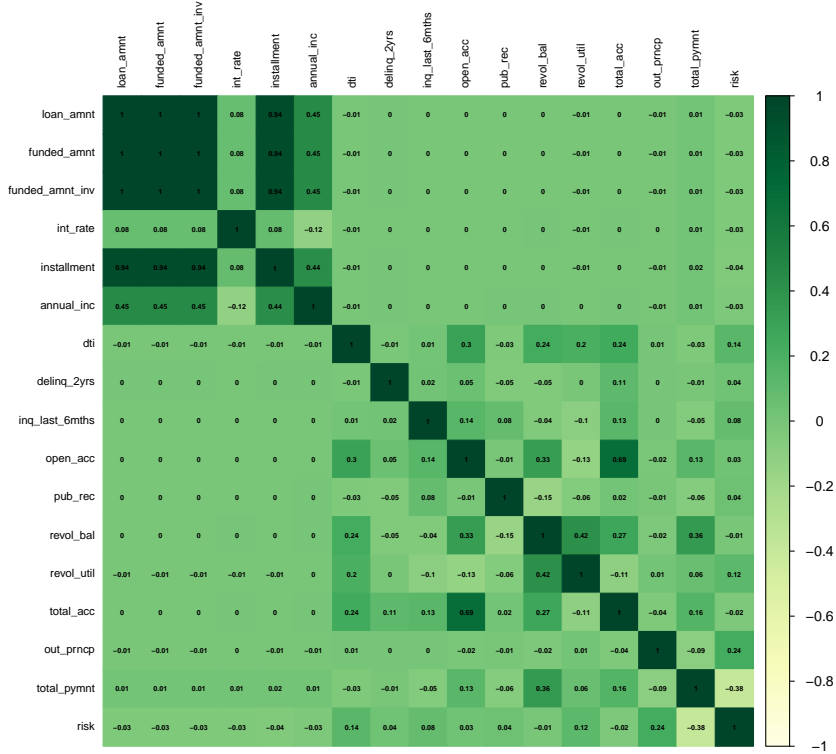


Figure 3: Correlation Heatmap

Also, features like *funded_amnt*, *funded_amnt_inv*, and *installment* were identified as highly correlated and were subsequently dropped from both training and test datasets to reduce multicollinearity, which could skew model results.

The Employment Length Conversion (*emp_length*), was where the unique values of the *emp_length* feature were inspected, then a mapping was defined to convert these categorical text values into numerical format, representing the number of years in employment, this conversion was applied to both training and test datasets. After the mapping, this feature was explicitly cast to numeric to ensure proper format. Furthermore, the loan term conversion (*term*) was when the feature was converted from a text format (e.g., "36 months") into a pure numeric format representing the number of months, this was achieved using the *gsub* function to remove non-numeric characters and then converting the result to numeric.

Some categorical features were transformed into numerical format using factorization, like: Grade, Address State, Home Ownership, Verification Status, Issue Date, and Purpose, then these features were converted into numeric by first turning them into factors and then mapping these factors to their underlying integer codes, this process was applied to both training and test datasets.

To deal with the *earliest_cr_line* feature which consists of time data, we decided to convert each observation into the difference in days between the current month (January 2024) and the observation itself, obtaining in this way a numerical variable.

A new binary feature was created to indicate whether an individual has been employed for 5 or more years, which could imply employment stability "Employment Stability" (*employment_stability*), then a new feature was calculated by dividing the annual income by the loan amount for each individual in both training and test datasets: "Income to Loan Ratio" (*income_to_loan_ratio*) which represents the proportion of income to the loan amount, which might be indicative of the ability to repay the loan.

A further cleaning process was needed in order to reduce dimensionality and focus on the most impactful variables by dropping high dispersion categorical features, these features with high dispersion or those considered less relevant, such as *emp_title* and *loan_status* were dropped from the training dataset. *emp_title* was also removed from the test dataset. The final structure display was when the structure of both the training and test datasets were displayed by using the *str()* function, providing an overview of the final dataset ready for modeling.

4 Modelling

4.1 Feature Selection

For the feature selection, we started by identifying and selecting numerical columns in the train dataset, creating a new variable called *numerical_columns*.

Then we calculated the correlations of each numerical feature with the target

variable *risk* and sorted their absolute values in descending order. This helps us understand which features have a stronger linear relationship with the target (*risk*), which is excluded from this sorted list (Figure 4).

After that, we created a horizontal bar plot to visually represent these correlations, aiding in the ulterior feature selection.

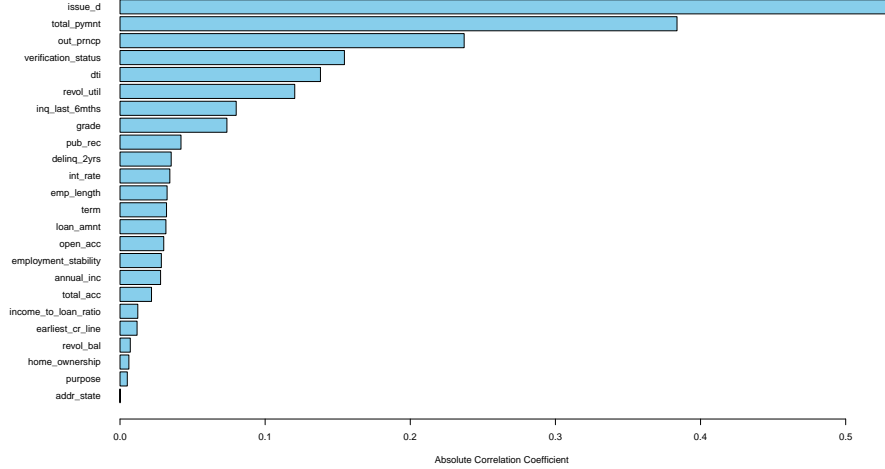


Figure 4: Correlation with the Target Feature 'risk'

We can see the highest correlations with the target variable are with the variables *total_pymnt*, *out_prncp*, *dti*, *revol_util*, and *inq_last_6mths*.

Then we proceeded to remove the weakly correlated features. We considered that 0.001 was the threshold, and only one variable (*addr_state*) with a correlation below that threshold was dropped from both training and testing datasets.

Then we dropped other variables due to different reasons:

- *emp_title*: due to high dispersion;
- *loan_status*: since it represents our target feature *risk*.

4.2 Logistic Regression

We started the modeling part by setting the seed to 123 to ensure the reproducibility of the results and splitting the data into training and test sets. A partition is created to ensure that two-thirds of the data is used for training.

Then we fit a logistic regression model to the training data with *risk* as the response variable and all other variables as predictors. Then we used a glm function with a binomial family and logit link function, suitable for binary classification.

Then we displayed a summary of the logistic regression model, which includes coefficient estimates, standard errors, z-values, and p-values, among other statistics.

Subsequently, we used the model to predict the risk on the test set. Predictions are probabilities, which are then converted into binary predictions using a threshold of 0.5. Then the performance of the model is evaluated using accuracy, calculated from a confusion matrix. A classification report is also printed, which includes metrics like Sensitivity (or Recall), Specificity, Positive Predictive Value (or Precision) and Negative Predictive Value, among others.

Additionally, we generated a confusion matrix and then normalized it to show proportions instead of absolute numbers, which helps in understanding the model's performance in differentiating between the classes.

	Reference	
Prediction	Non-Default	Default
Non-Default	0.8366906	0.1633094
Default	0.2617874	0.7382126

Table 3: Confusion Matrix Logistic Regression

Indicator	Values
Accuracy Rate	0.8105
Sensitivity	0.8980
Specificity	0.6213
Pos Pred Value	0.8367
Neg Pred Value	0.7382
Balanced Accuracy	0.7597

Table 4: Performance measures Logistic Regression

Finally, a Receiver Operating Characteristic (ROC) curve is plotted, and the Area Under Curve (AUC) score is calculated. The ROC curve is a plot of the true positive rate against the false positive rate at various threshold settings, and the AUC score is a measure of the model's ability to distinguish between the classes.

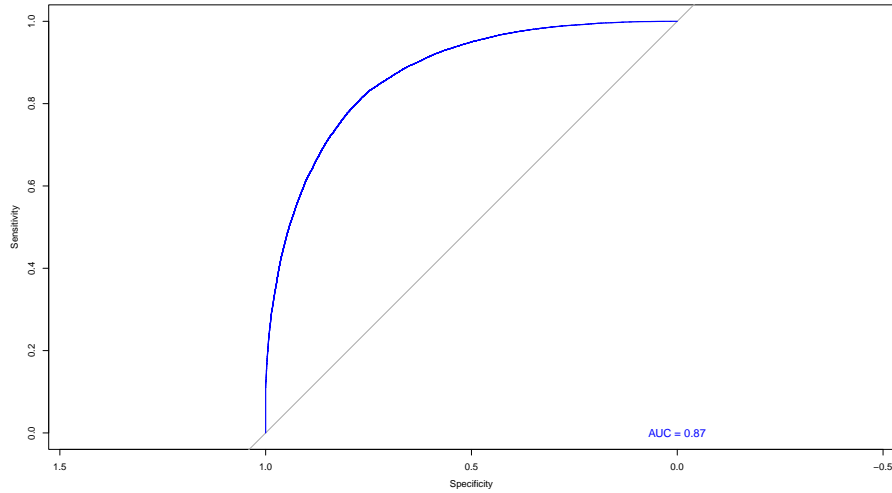


Figure 5: Logistic Regression ROC Curve

4.3 Random Forest

For the random forest model, we created one hundred trees and the result that we got for the Area Under the Curve was 0.9197253, which is higher than the 0.8731501 that we got on the logistic regression. This indicates that the random forest model should present better results.

However, when we are comparing the two models, we need to take into account other measures. Comparing the Random Forest model with the Logistic Regression, we can clearly see that the accuracy of the Random Forest model is higher, another value that is also higher is the specificity which is a very relevant measure for the problem we are solving, it allows us to correctly identify true negatives which will prevent big losses.

For this specific problem we can clearly see that we will have better results if we implement the Random Forest model instead of the Logistic Regression model.

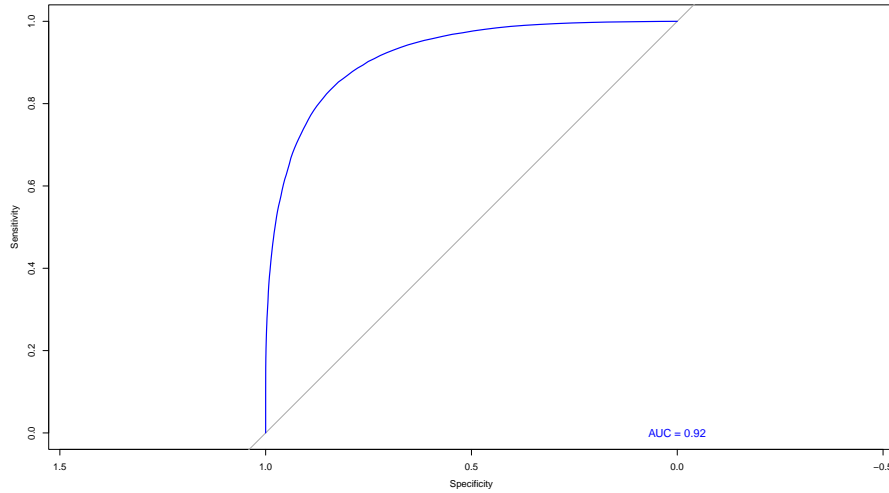


Figure 6: Random Forest ROC Curve

	Reference	
Prediction	Non-Default	Default
Non-Default	0.8552732	0.1447268
Default	0.1582959	0.8417041

Table 5: Confusion Matrix Random Forest

Indicator	Values
Accuracy Rate	0.8519
Sensitivity	0.9430
Specificity	0.6552
Pos Pred Value	0.8553
Neg Pred Value	0.8417
Balanced Accuracy	0.7978

Table 6: Performance measures Random Forest

4.4 Deep Neural Network

Finally, the Neural Network model when compared to the Random Forest model, does not outperform it in any measure, like the sensitivity and the Area Under the Curve and it underperformed in the specificity with 0.2568 when compared with 0.6552 of the Random Forrest.

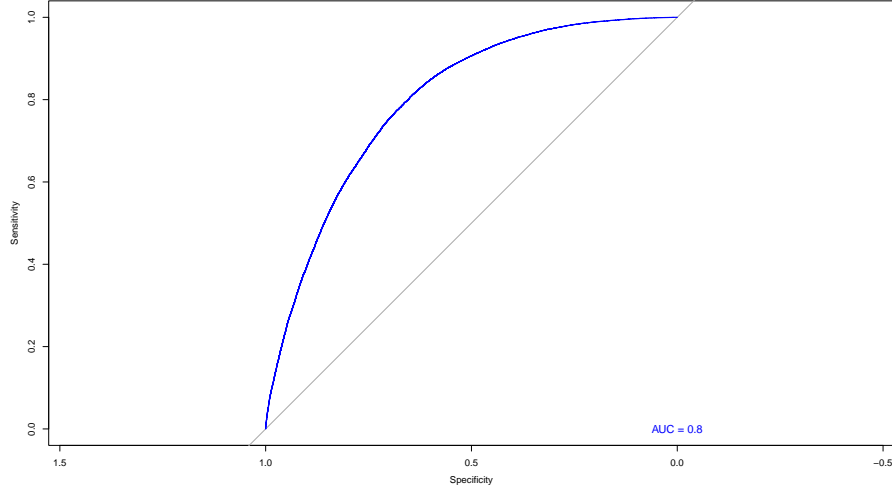


Figure 7: Deep Neural Network ROC Curve

	Reference	
Prediction	Non-Default	Default
Non-Default	0.7339719	0.2660281
Default	0.3096887	0.6903113

Table 7: Confusion Matrix Deep Neural Network

Indicator	Values
Accuracy Rate	0.7288
Sensitivity	0.9468
Specificity	0.2568
Pos Pred Value	0.7340
Neg Pred Value	0.6903
Balanced Accuracy	0.6018

Table 8: Performance measures Deep Neural Network

Comparing the three models that we tested the one with the worst results was the Neural Network, with such low specificity if we applied this model we would incur in very big losses because we would not be able to accurately predict the True Negatives. The Logistic Regression managed to outperform the Neural Network but is outperformed by the Random Forrest which has a better AUC and most importantly the specificity.

In order to have better results we choose to implement the Random Forrest model.

5 Conclusion

In this analysis, we employed three distinct Machine Learning Models (Logistic Regression, Random Forest, and a Deep Neural Network) to predict the probability of default in a credit dataset. The Logistic Regression model demonstrated a commendable overall accuracy of 81.05%, with a balanced accuracy of 75.97%. It exhibited strong sensitivity (true positive rate) of 89.80%, emphasizing its effectiveness in correctly identifying instances of default. The Random Forest model outperformed the Logistic Regression with an impressive accuracy of 85.1% and a balanced accuracy of 79.78%. It showcased high sensitivity (94.25%), indicating its proficiency in capturing true positives. However, the Deep Neural Network, while achieving an overall accuracy of 72.88%, struggled with specificity, resulting in a lower balanced accuracy of 60.18%. This suggests that the deep neural network may be more prone to false positives. In summary, the random forest model emerged as the most robust performer in predicting default probabilities, striking a balance between accuracy and sensitivity. These findings underscore the importance of selecting an appropriate model based on the specific needs and priorities of the credit risk assessment task. Future work could explore model fine-tuning and ensemble techniques to further enhance predictive performance and provide more robust risk assessment in the realm of credit evaluation.

References

- [1] Adams, K., & Shin, D. (2022), *The history of credit score algorithms and how they became the lender standard*.
Available at: <https://www.marketplace.org/shows/marketplace-tech/the-history-of-credit-score-algorithms-and-how-they-became-the-lender-standard/>
- [2] CFI Team (2022), *Credit Risk Analysis Models, Tools used to determine the probability of default of a potential borrower*.
Available at: <https://corporatefinanceinstitute.com/resources/commercial-lending/credit-risk-analysis-models/#:~:text=Credit%20risk%20arises%20when%20a,debt%20extended%20to%20a%20borrower.>
- [3] Jain, J. & Banka, P. (2022), *Dealing with outliers and missing values in a dataset*.
Available at: <https://www.neenopal.com/dealing-with-outliers-and-missing-values-in-a-dataset.html>
- [4] Mongelard, J. (2019), *A brief history in credit risk in banking*.
Available at: <https://www.icaew.com/technical/financial-services/financial-services-faculty/fs-focus-magazine/previous-edition-s-of-fs-focus/fs-focus-2019-issues/november-2019/brief-history-in-credit-risk-in-banking>
- [5] Pilz, J., Melas, V.B., & Bathke, A. (2023). *Statistical Modeling and Simulation for Experimental Design and Machine Learning Applications*. Springer Nature, pp.245-265.
- [6] Zach (2022), *How to Perform Exploratory Data Analysis in R (With Example)*.
Available at: <https://www.statology.org/exploratory-data-analysis-in-r/>

Code

```

1  if (!require("psych")) install.packages("psych")
2  library(psych)
3  if (!require("dplyr")) install.packages("dplyr")
4  library(dplyr)
5  if (!require("gmodels")) install.packages("gmodels")
6  library(gmodels)
7  if (!require("tidyverse")) install.packages("tidyverse")
8  library(tidyverse)
9  if (!require("corrplot")) install.packages("corrplot")
10 library(corrplot)
11 if (!require("caret")) install.packages("caret")
12 library(caret)
13 if (!require("pROC")) install.packages("pROC")
14 library(pROC)
15 if (!require("ROCR")) install.packages("ROCR")
16 library(ROCR)
17 if (!require("randomForest")) install.packages("randomForest")
18 library(randomForest)
19 if (!require("keras")) install.packages("keras")
20 library(keras)
21 if (!require("tensorflow")) install.packages("tensorflow")
22 library(tensorflow)
23 if (!require("lubridate")) install.packages("lubridate")
24 library(lubridate)
25
26
27 # Loading Data
28 # -----
29 train_dataset <- read.csv("/Users/riccardogurzu/Desktop/predictive
↳ analytics/train_validation_kaggle.csv",
30                           na.strings = c("", "n/a"),
31                           stringsAsFactors = T)
32 # Remove the 'id' column from the dataset
33 train_dataset <- train_dataset[, !(names(train_dataset) %in% c("id"))]
34
35 test_dataset <- read.csv("/Users/riccardogurzu/Desktop/predictive
↳ analytics/unseen_kaggle.csv",
36                          na.strings = c("", "n/a"),
37                          stringsAsFactors = T)
38
39 # View datasets
40 View(train_dataset)
41 View(test_dataset)
42
43 #Check details
44 summary(train_dataset)
45 summary(test_dataset)
46 head(train_dataset, n=5)
47 tail(train_dataset, n=5)
48 describe(train_dataset)
49 describe(test_dataset)
50
51 str(train_dataset)
52 str(test_dataset)
53

```

```

54
55 # Create a bar plot of the 'risk' to check the balance of the dataset
56 barplot(table(train_dataset$risk),
57         col = c("lightblue"),
58         main = "Risk",
59         ylab = "Frequency")
60
61
62 # -----
63 #Data #Pre-Processing
64 # -----
65 ##Data-Type
66 str(train_dataset)
67 str(test_dataset)
68
69 # -----
70 ##Missing Values
71
72 ###Train Dataset
73 sapply(train_dataset, function(x) sum(is.na(x)))
74
75 ###emp_title - 29563 missing values
76 train_dataset$emp_title = factor(train_dataset$emp_title,
77                                 levels=c(levels(train_dataset$emp_title),
78                                           ↪ 'Unemployed'))
79
80 ###emp_length - 22615 missing values
81 train_dataset$emp_length[is.na(train_dataset$emp_length)] <- '< 1 year'
82
83 ###dti - 148 missing values
84 mean(train_dataset$dti, na.rm = TRUE)
85 train_dataset$dti[is.na(train_dataset$dti)] <- mean(train_dataset$dti, na.rm =
86 ↪ TRUE)
87
88 ###inq_last_6mths - 1 missing values
89 mean(train_dataset$inq_last_6mths, na.rm = TRUE)
90 train_dataset$inq_last_6mths[is.na(train_dataset$inq_last_6mths)] <-
91 ↪ mean(train_dataset$inq_last_6mths, na.rm = TRUE)
92
93 ###revol_util - 213 missing values
94 mean(train_dataset$revol_util, na.rm = TRUE)
95 train_dataset$revol_util[is.na(train_dataset$revol_util)] <-
96 ↪ mean(train_dataset$revol_util, na.rm = TRUE)
97
98 sapply(train_dataset, function(x) sum(is.na(x)))
99
100
101 ###Test Dataset
102 str(test_dataset)
103 sapply(test_dataset, function(x) sum(is.na(x)))
104
105 ### dropping the row where 'loan_amnt' has missing values since 180 rows are
106 ↪ not correctly separated
107 ### due to values in emp_title with double quotes which contain commas

```

```

106 test_dataset <- test_dataset[complete.cases(test_dataset$loan_amnt), ]
107
108 sapply(test_dataset, function(x) sum(is.na(x)))
109
110 ###emp_title - 1224 missing values
111 test_dataset$emp_title = factor(test_dataset$emp_title,
112                                levels=c(levels(test_dataset$emp_title),
113                                          ↪ 'Unemployed'))
113 test_dataset$emp_title[is.na(test_dataset$emp_title)] <- 'Unemployed'
114
115 ###emp_length - 1210 missing values
116 test_dataset$emp_length[is.na(test_dataset$emp_length)] <- '< 1 year'
117
118 ###dti - 18 missing values
119 test_dataset$dti <- as.numeric(test_dataset$dti)
120 mean(test_dataset$dti, na.rm = TRUE)
121 test_dataset$dti[is.na(test_dataset$dti)] <- mean(test_dataset$dti, na.rm =
122 ↪ TRUE)
123
124 ###revol_util - 17 missing values
125 test_dataset$revol_util <- as.numeric(test_dataset$revol_util)
126 mean(test_dataset$revol_util, na.rm = TRUE)
127 test_dataset$revol_util[is.na(test_dataset$revol_util)] <-
128 ↪ mean(test_dataset$revol_util, na.rm = TRUE)
129
130 sapply(test_dataset, function(x) sum(is.na(x)))
131
132 # -----
133
134 ##Outliers
135 boxplot(train_dataset)
136 str(train_dataset)
137
138 numerical_columns <- sapply(train_dataset, is.numeric)
139
140 #print the boxplot for the numerical features
141 for (col in names(train_dataset[, numerical_columns])) {
142   boxplot(train_dataset[, col],
143           main = col,
144           col = "skyblue",
145           border = "black")
146 }
147
148 #print the histogram for the numerical features
149 for (col in names(train_dataset[, numerical_columns])) {
150   hist(train_dataset[, col],
151        main = col,
152        col = "skyblue",
153        border = "black",
154        xlim = c(min(train_dataset[, col]), max(train_dataset[, col])),
155        breaks = 30)
156 }
157
158 # for categorical features
159 categorical_columns <- sapply(train_dataset, is.factor)

```

```

160 # exclude emp_title, issue_d and earliest_cr_line due to the high number of
    ↪ levels
161 exclude_features <- c("emp_title", "issue_d", "earliest_cr_line")
162
163 # Identify the indices of the categorical columns to keep
164 categorical_columns <- !(names(train_dataset) %in% exclude_features) &
    ↪ categorical_columns
165
166 # Loop through each categorical column and plot bar charts
167 for (col in names(train_dataset[, categorical_columns])) {
168   barplot(table(train_dataset[, col]),
169     main = col,
170     col = "skyblue",
171     xlab = "Categories",
172     ylab = "Count")
173 }
174
175
176 # Define threshold values
177 threshold_values <- list(
178   int_rate = 25,
179   annual_inc = 0.3 * 10^6,
180   dti = 50,
181   delinq_2yrs = 5,
182   inq_last_6mths = 3,
183   open_acc = 30,
184   pub_rec = 3,
185   revol_bal = 60000,
186   revol_util = 110,
187   total_acc = 70,
188   out_prncp = 10000,
189   total_pymnt = 50000,
190   credit_line_duration = 20000,
191   income_to_loan_ratio = 40
192 )
193
194 # Loop through each column and apply the corresponding threshold for the Train
    ↪ Dataset
195 for (column in names(train_dataset)) {
196   if (column %in% names(threshold_values)) {
197     train_dataset <- train_dataset[train_dataset[, column] <=
    ↪ threshold_values[[column]], ]
198   }
199 }
200
201 # Loop through each column and apply the corresponding threshold for the Test
    ↪ Dataset
202 for (column in names(test_dataset)) {
203   if (column %in% names(threshold_values)) {
204     test_dataset <- test_dataset[test_dataset[, column] <=
    ↪ threshold_values[[column]], ]
205   }
206 }
207
208 #print the boxplot for the numerical features after dropped the outliers
209 for (col in names(train_dataset[, numerical_columns])) {
210   boxplot(train_dataset[, col],

```

```

211     main = col,
212     col = "skyblue",
213     border = "black")
214 }
215
216 #print the histogram for the numerical features after dropped the outliers
217 for (col in names(train_dataset[, numerical_columns])) {
218     hist(train_dataset[, col],
219         main = col,
220         col = "skyblue", border = "black",
221         xlim = c(min(train_dataset[, col]), max(train_dataset[, col])),
222         breaks = 30)
223 }
224
225 # -----
226 #Feature Engineering
227 # -----
228
229 # Select the columns with numerical data
230 numerical_columns <- sapply(train_dataset, is.numeric)
231
232 # Create a correlation matrix
233 correlation_matrix <- cor(train_dataset[, numerical_columns])
234 correlation_matrix
235
236 # Plot the correlation matrix heatmap
237 corrplot(correlation_matrix,
238     method = "color",
239     tl.col = "black",
240     col = COL1('YlGn'),
241     addCoef.col = "black",
242     tl.cex = 0.7,
243     number.cex = 0.4)
244
245 #drop high correlated features (>|0.7|)
246 train_dataset <- train_dataset[, -which(names(train_dataset) %in%
247     ↪ c("funded_amnt", "funded_amnt_inv", "installment"))]
248 test_dataset <- test_dataset[, -which(names(test_dataset) %in%
249     ↪ c("funded_amnt", "funded_amnt_inv", "installment"))]
250
251 ### Feature Transformation
252 str(train_dataset)
253 str(test_dataset)
254
255 # removing the dot in total_pymnt. to match train dataset
256 colnames(test_dataset)[colnames(test_dataset) == 'total_pymnt.'] <-
257     ↪ 'total_pymnt'
258 # Convert total_pymnt into numeric
259 test_dataset$total_pymnt <- as.numeric(test_dataset$total_pymnt)
260
261 # Convert emp_length into numeric
262 unique(train_dataset$emp_length)
263 mapping <- c(
264     "< 1 year" = 0,
265     "1 year" = 1,

```

```

265   "2 years" = 2,
266   "3 years" = 3,
267   "4 years" = 4,
268   "5 years" = 5,
269   "6 years" = 6,
270   "7 years" = 7,
271   "8 years" = 8,
272   "9 years" = 9,
273   "10+ years" = 10
274 )
275
276 train_dataset <- train_dataset %>%
277   mutate(emp_length = case_when(
278     emp_length %in% names(mapping) ~ mapping[as.character(emp_length)],
279     TRUE ~ NA_real_
280   ))
281 train_dataset$emp_length <- as.numeric(train_dataset$emp_length)
282 unique(train_dataset$emp_length)
283
284 test_dataset <- test_dataset %>%
285   mutate(emp_length = case_when(
286     emp_length %in% names(mapping) ~ mapping[as.character(emp_length)],
287     TRUE ~ NA_real_
288   ))
289 test_dataset$emp_length <- as.numeric(test_dataset$emp_length)
290 unique(test_dataset$emp_length)
291
292 # Convert term into numeric
293 train_dataset$term <- as.numeric(gsub("[^0-9]+", "", train_dataset$term))
294 test_dataset$term <- as.numeric(gsub("[^0-9]+", "", test_dataset$term))
295
296 # Transform the levels for the categoric features (except emp_title and
297   ↪ earliest_cr_line) into numeric
298 unique(train_dataset$grade)
299 train_dataset$grade <- as.numeric(factor(train_dataset$grade,
300                                         levels =
301                                           ↪ levels(train_dataset$grade)))
302
303 unique(test_dataset$grade)
304 test_dataset$grade <- as.numeric(factor(test_dataset$grade,
305                                         levels = levels(test_dataset$grade)))
306
307
308 unique(train_dataset$addr_state)
309 train_dataset$addr_state <- as.numeric(factor(train_dataset$addr_state,
310                                               levels =
311                                                 ↪ levels(train_dataset$addr_state)))
312
313 unique(test_dataset$addr_state)
314 test_dataset$addr_state <- as.numeric(factor(test_dataset$addr_state,
315                                               levels =
316                                                 ↪ levels(test_dataset$addr_state)))
317
318 unique(train_dataset$home_ownership)

```

```

318 train_dataset$home_ownership <-
↳ as.numeric(factor(train_dataset$home_ownership,
319 levels =
↳ levels(train_dataset$home_ownership)))
320 unique(test_dataset$home_ownership)
321 test_dataset$home_ownership <- as.numeric(factor(test_dataset$home_ownership,
322 levels =
↳ levels(test_dataset$home_ownership)))
323
324 unique(train_dataset$verification_status)
325 train_dataset$verification_status <-
326 ↳ as.numeric(factor(train_dataset$verification_status,
327 levels =
↳ levels(train_dataset$verification_status)))
328
329 unique(test_dataset$verification_status)
330 test_dataset$verification_status <-
331 ↳ as.numeric(factor(test_dataset$verification_status,
332 levels =
↳ levels(test_dataset$verification_status)))
333
334 unique(train_dataset$purpose)
335 train_dataset$purpose <- as.numeric(factor(train_dataset$purpose,
336 levels =
↳ levels(train_dataset$purpose)))
337
338 unique(test_dataset$purpose)
339 test_dataset$purpose <- as.numeric(factor(test_dataset$purpose,
340 levels =
↳ levels(test_dataset$purpose)))
341
342 # Create 'employment_stability' feature
343 train_dataset$employment_stability <- ifelse(train_dataset$emp_length >= 5, 1,
↳ 0)
344 test_dataset$employment_stability <- ifelse(test_dataset$emp_length >= 5, 1,
↳ 0)
345
346 # Create 'income_to_loan_ratio' feature
347 train_dataset$income_to_loan_ratio <-
↳ train_dataset$annual_inc/train_dataset$loan_amnt
348 test_dataset$annual_inc <- as.numeric(test_dataset$annual_inc)
349 test_dataset$income_to_loan_ratio <-
↳ test_dataset$annual_inc/test_dataset$loan_amnt
350
351
352 ## To convert earliest_cr_line into numeric, we modify the feature computing
353 ↳ the diff with the current month
354 # Convert the 'earliest_cr_line' column to a Date type with custom format
355 train_dataset$earliest_cr_line <- dmy(paste0("01-",
↳ train_dataset$earliest_cr_line))
356 unique(train_dataset$earliest_cr_line)
357 # Adjust two-digit years using if_else
358 train_dataset$earliest_cr_line <- if_else(year(train_dataset$earliest_cr_line)
↳ <= 2024,
train_dataset$earliest_cr_line,

```

```

359                                     train_dataset$earliest_cr_line -
                                     ↪ years(100))
360 unique(train_dataset$earliest_cr_line)
361
362 test_dataset$earliest_cr_line <- as.Date(test_dataset$earliest_cr_line, origin
                                     ↪ = "1899-12-30")
363 unique(train_dataset$earliest_cr_line)
364 unique(test_dataset$earliest_cr_line)
365
366 # Calculate the difference in months between each date and the reference date
367 current_date <- dmy("01-01-2024")
368 train_dataset$earliest_cr_line <-
                                     ↪ as.numeric(interval(train_dataset$earliest_cr_line, current_date) /
                                     ↪ months(1))
369 test_dataset$earliest_cr_line <-
                                     ↪ as.numeric(interval(test_dataset$earliest_cr_line, current_date) /
                                     ↪ months(1))
370 unique(train_dataset$earliest_cr_line)
371 unique(test_dataset$earliest_cr_line)
372
373 unique(train_dataset$issue_d)
374 unique(test_dataset$issue_d)
375
376 # Convert numerical labels to Date type
377 test_dataset$issue_d <- as.Date(test_dataset$issue_d, origin = "1899-12-30")
378 unique(test_dataset$issue_d)
379 # since levels of issue_d between train and test are different, we cannot use
                                     ↪ label encoding
380 # thus, we modify the issue_d feature into the diff in months between issue_d
                                     ↪ and current month
381
382 # Convert 'issue_d' to a Date type
383 train_dataset$issue_d <- dmy(paste0("01-", train_dataset$issue_d))
384 unique(train_dataset$issue_d)
385
386 current_date <- dmy("01-01-2024")
387
388 # Calculate the difference in months between each date and the current date
389 train_dataset$issue_d <- as.numeric(interval(train_dataset$issue_d,
                                     ↪ current_date) / months(1))
390 test_dataset$issue_d <- as.numeric(interval(test_dataset$issue_d,
                                     ↪ current_date) / months(1))
391 unique(train_dataset$issue_d)
392 unique(test_dataset$issue_d)
393
394
395 # Display the result
396 str(train_dataset)
397 str(test_dataset)
398
399
400 ### Plot the correlation between the target feature 'Risk' and the other
                                     ↪ features
401
402 # Select the columns with numerical data
403 numerical_columns <- sapply(train_dataset, is.numeric)
404

```



```

405 # Calculate the correlation with the target variable
406 correlation_with_target <- sapply(train_dataset[, numerical_columns],
407                                   function(x) cor(x, train_dataset$risk))
408
409 # Order correlations by absolute values in descending order
410 sorted_correlations <- sort(abs(correlation_with_target), decreasing = FALSE)
411
412 # Exclude the 'risk' variable
413 sorted_correlations <- sorted_correlations[!(names(sorted_correlations) %in%
414   ↪ "risk")]
415
416 # Set larger plot margins to fit the y labels
417 par(mar = c(5, 10, 2, 2))
418
419 # Create a barplot with the correlations to 'risk'
420 barplot(sorted_correlations, horiz = TRUE, col = "skyblue",
421         xlab = "Absolute Correlation Coefficient",
422         las = 1)
423
424 # Set plot margins to default
425 par(mar = c(5, 4, 4, 2))
426
427 # Extract column names with correlation less than 0.01 and drop them
428 columns_to_drop <- names(correlation_with_target[abs(correlation_with_target)
429   ↪ < 0.001]) #only addr_state meet this condition
430 train_dataset <- train_dataset[, !names(train_dataset) %in% columns_to_drop]
431 test_dataset <- test_dataset[, !names(test_dataset) %in% columns_to_drop]
432
433 #drop emp_title due to high dispersion
434 #drop loan_status since represent our target feature 'risk'
435 #drop issue_d since if we include it we get for some reason a perfect
436   ↪ performance
437 train_df <- train_dataset[, -which(names(train_dataset) %in% c("emp_title",
438   ↪ "loan_status", "issue_d"))]
439 test_df <- test_dataset[, -which(names(test_dataset) %in% c("emp_title",
440   ↪ "issue_d"))]
441
442 # Display the result
443 str(train_df)
444 str(test_df)
445
446 # -----
447 # Logistic regression
448 # -----
449
450 #Breaking Data into Training and Test Sample
451 set.seed(123)
452
453 # Data splitting
454 index <- createDataPartition(train_df$risk, p = 2/3, list = FALSE)
455 train <- train_df[index, ]
456 test <- train_df[-index, ]
457
458 # Fit the logistic regression model
459 model_log <- glm(risk ~ .,
460                 data = train,

```

```

457         family = binomial(link = "logit"))
458
459     # Summary of the model
460     summary(model_log)
461
462     # Make predictions on the test set
463     y_pred <- predict(model_log, newdata = test, type = "response")
464
465     # Convert probabilities to binary predictions (assuming a threshold of 0.5)
466     y_pred_binary <- ifelse(y_pred > 0.5, 1, 0)
467     y_pred_binary <- as.factor(y_pred_binary)
468
469     ## Evaluate the model
470     test$risk <- as.factor(test$risk)
471     # Accuracy
472     accuracy <- confusionMatrix(data = y_pred_binary,
473                                reference = test$risk)$overall["Accuracy"]
474     cat("Accuracy:", accuracy, "\n")
475
476     # Classification Report
477     conf_matrix <- confusionMatrix(data = y_pred_binary, reference = test$risk)
478     print(conf_matrix)
479
480     # Confusion Matrix
481     conf_matrix <- as.table(conf_matrix)
482     conf_matrix <- prop.table(conf_matrix, 1)
483
484     # Plot the correlation matrix
485     conf_matrix
486
487     #Plot ROC Curve
488     roccurve <- roc(test$risk ~ y_pred)
489     auc_score <- auc(roccurve)
490     cat("Logistic Regression AUC Score:", auc_score, "\n")
491
492     plot(roccurve, col='blue')
493     # Add AUC value as text annotation
494     text(0, 0, paste("AUC =", round(auc_score, 2)), col="blue", cex=1.2)
495
496
497
498     # -----
499     # Random Forest
500     # -----
501
502     #Breaking Data into Training and Test Sample
503     set.seed(123)
504
505     # Data splitting
506     index <- createDataPartition(train_df$risk, p = 2/3, list = FALSE)
507     train <- train_df[index, ]
508     test <- train_df[-index, ]
509
510     train$risk <- as.factor(train$risk)
511
512     # Fit the Random Forest model
513     model_rf <- randomForest(risk ~ .,

```

```

514         data = train,
515         ntree = 100)
516
517 # Make predictions on the test set
518 y_pred_rf <- predict(model_rf, type = "prob", test)
519
520 # Convert probabilities to binary predictions (assuming a threshold of 0.5)
521 y_pred_binary_rf <- ifelse(y_pred_rf[,2] > 0.5, 1, 0)
522 y_pred_binary_rf <- as.factor(y_pred_binary_rf)
523
524 ## Evaluate the Random Forest model
525 test$risk <- as.factor(test$risk)
526
527 # Accuracy
528 accuracy_rf <- confusionMatrix(data = y_pred_binary_rf, reference =
529   ↪ test$risk)$overall["Accuracy"]
529 cat("Random Forest Accuracy:", accuracy_rf, "\n")
530
531 # Classification Report
532 conf_matrix_rf <- confusionMatrix(data = y_pred_binary_rf, reference =
533   ↪ test$risk)
533 print(conf_matrix_rf)
534
535 # Confusion Matrix
536 conf_matrix_rf <- as.table(conf_matrix_rf)
537 conf_matrix_rf <- prop.table(conf_matrix_rf, 1)
538
539 # Plot the confusion matrix
540 conf_matrix_rf
541
542 # Plot ROC Curve for Random Forest
543 roccurve_rf <- roc(test$risk ~ y_pred_rf[,2])
544 # AUC Score for Random Forest
545 auc_score_rf <- auc(roccurve_rf)
546 cat("Random Forest AUC Score:", auc_score_rf, "\n")
547 plot(roccurve_rf, col = 'blue')
548 # Add AUC value as text annotation
549 text(0, 0, paste("AUC =", round(auc_score_rf, 2)), col="blue", cex=1.2)
550
551
552 # -----
553 # Deep Neural Network
554 # -----
555
556 install_tensorflow(envname = "r-tensorflow")
557
558 tf$constant("Hello TensorFlow!")
559
560 #Breaking Data into Training and Test Sample
561 set.seed(123)
562
563 # Data splitting
564 index <- createDataPartition(train_df$risk, p = 2/3, list = FALSE)
565 train <- train_df[index, ]
566 test <- train_df[-index, ]
567
568 y_train <- as.numeric(train$risk)

```

```

569 y_test <- as.numeric(test$risk)
570 x_train <- train %>% select(-risk)
571 x_test <- test %>% select(-risk)
572
573 # Normalize hte input features
574 x_train <- as.matrix(apply(x_train, 2, function(x) (x-min(x))/(max(x) -
↪ min(x))))
575 x_test <- as.matrix(apply(x_test, 2, function(x) (x-min(x))/(max(x) -
↪ min(x))))
576
577
578 # Fit the Neural Network model
579 model_nn <- keras_model_sequential() %>%
580   layer_dense(units = 64, activation = "relu", input_shape = ncol(train) - 1)
↪ %>%
581   layer_dense(units = 32, activation = "relu") %>%
582   layer_dense(units = 1, activation = "sigmoid")
583
584 compile(model_nn, optimizer = "adam",
585         loss = "binary_crossentropy",
586         metrics = c("accuracy"))
587
588 # Train the model
589 history <- fit(
590   model_nn,
591   x = x_train,
592   y = y_train,
593   shuffle = T,
594   epochs = 10,
595   batch_size = 32,
596   validation_split = 0.2
597 )
598
599 # Make predictions on the test set
600 y_pred_prob_nn <- predict(model_nn, x_test)
601
602 # Convert probabilities to binary predictions
603 y_pred_binary_nn <- ifelse(y_pred_prob_nn > 0.5, 1, 0)
604 y_pred_binary_nn <- as.factor(y_pred_binary_nn)
605
606 # Evaluate the neural network model
607 y_test <- as.factor(y_test)
608 # Accuracy
609 accuracy_nn <- confusionMatrix(data = y_pred_binary_nn, reference =
↪ y_test)$overall["Accuracy"]
610 cat("Neural Network Accuracy:", accuracy_nn, "\n")
611
612 # Classification Report
613 conf_matrix_nn <- confusionMatrix(data = y_pred_binary_nn, reference = y_test)
614 print(conf_matrix_nn)
615
616 # Confusion Matrix
617 conf_matrix_nn <- as.table(conf_matrix_nn)
618 conf_matrix_nn <- prop.table(conf_matrix_nn, 1)
619
620 # Plot the confusion matrix
621 conf_matrix_nn

```

```

622 # Plot ROC Curve for Neural Network
623 roccurve_nn <- roc(y_test, as.numeric(y_pred_prob_nn))
624 # AUC Score for Neural Network
625 auc_score_nn <- auc(roccurve_nn)
626 cat("Neural Network AUC Score:", auc_score_nn, "\n")
627 plot(roccurve_nn, col = 'blue')
628 # Add AUC value as text annotation
629 text(0, 0, paste("AUC =", round(auc_score_nn, 2)), col="blue", cex=1.2)
630
631
632
633 # -----
634 # Deployment
635 # -----
636 # Based on the evaluation metrics, Random Forest is the model which perform
637 ↪ better
638
639 # Make predictions on the unseen file using the trained model
640 unseen_df <- test_df[, !(names(test_df) %in% c("id"))]
641 y_pred_rf <- predict(model_rf, type = "prob", unseen_df)
642
643 # Convert probabilities to binary predictions (assuming a threshold of 0.5)
644 y_pred_binary_rf <- ifelse(y_pred_rf[,2] > 0.5, 1, 0)
645 y_pred_binary_rf <- as.factor(y_pred_binary_rf)
646
647 RandomForest <- data.frame(id=test_df$id, risk=y_pred_binary_rf)
648 RandomForest
649
650 write.csv(RandomForest, "Group5.csv", row.names = FALSE)

```