# Predicting customers' defaults

### ING Lion's Den credit risk modeling exercise
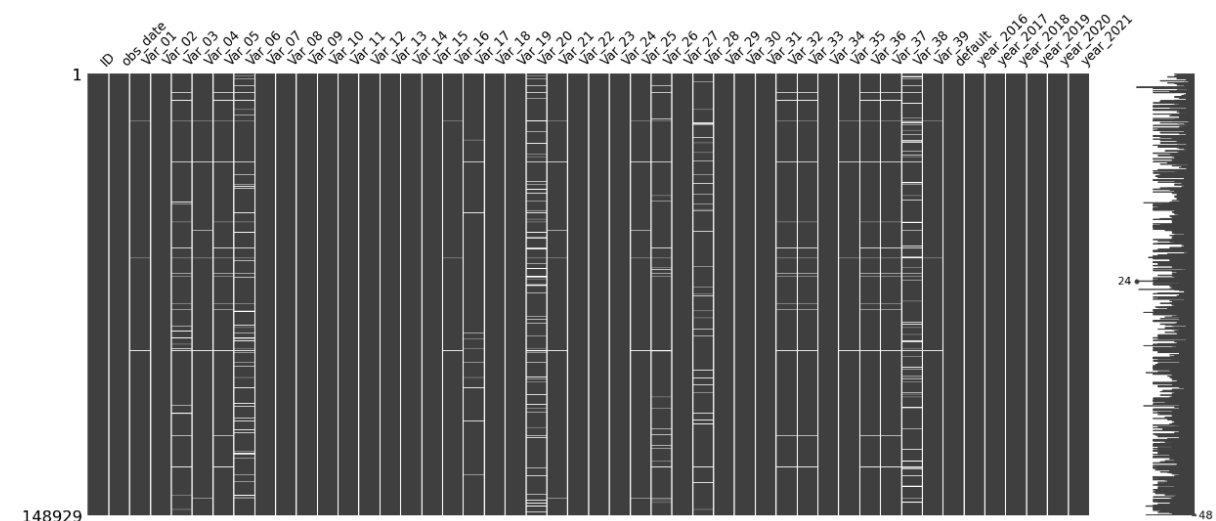### Monika Kaczan

## Introduction

The goal of this project is to predict defaults of Big Lion Bank's customers. I received annual data on borrowing companies including financial information and information whether a particular company defaulted within the next 12 months. Using this data, I build logit, ridge, KNN and XGBoost models. I compared their performance based on F1 score and other metrics and chose the best one.

## Data preprocessing

The raw data contained 148929 observations of unique customers observed in the years 2015-2021. It included 39 numerical variables describing their financial characteristics, time variable, and one dummy variable indicating performing company defaulting within the next 12 months.

<u>Handling missing values</u>



Visualization of missing values in the dataset and particular variables

Over 40,000 of observations (around 27%) contained at least one missing value. Therefore, I find simply deleting them as not reasonable. I decided to pursue following approach:
- deleting variables which had very high correlations with other variables AND had more than 1500 (~1% of entire dataset) NaNs: Var_03, Var_17, Var_20, Var_26, Var_38
- predicting missing values using linear regression for variables which had more than 1500 NaNs but only moderate correlations with other variables: Var_06 and Var_28

● filling NaNs with median values (as the distributions were skewed) for rest of variables which contained NaNs

Outliers detection and handling

Investigating variables' distributions showed that they contain observations which could be classified as outliers. This could be due to human errors or just the characteristics of such financial data. Nevertheless, outliers can significantly affect performance of many types of models. This is why I decided to cap and floor each column with 99th and 1st percentile respectively.

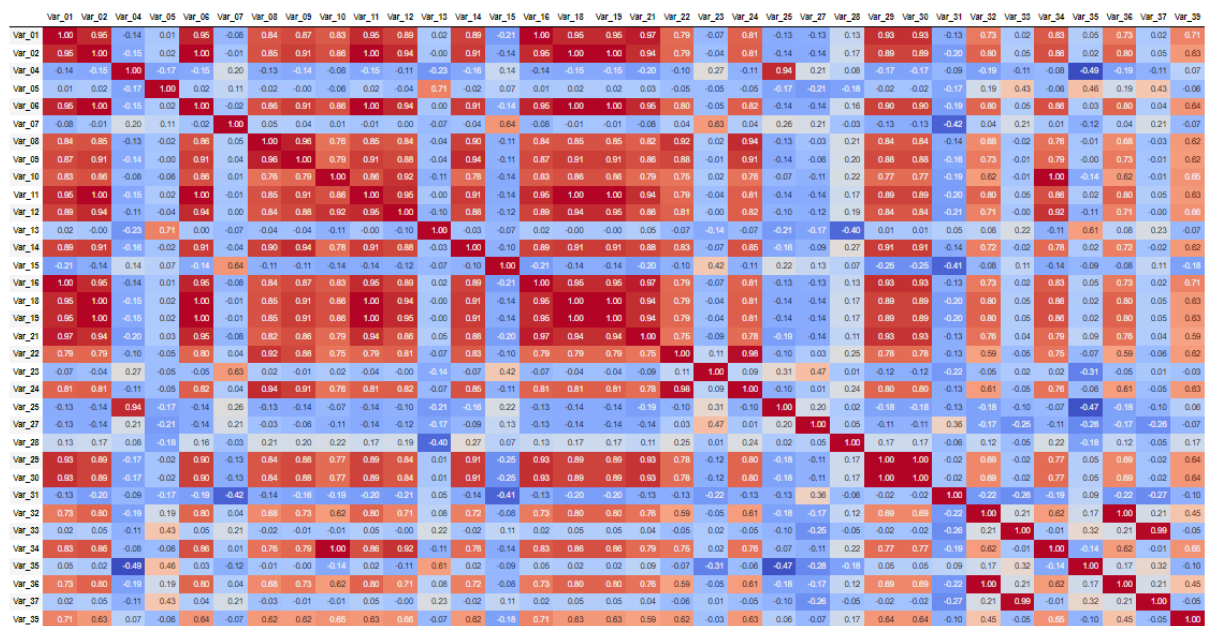Splitting data into training and testing samples

I split the data into training (70% of entire dataset) and testing (30% of entire dataset) samples. I made sure that distribution of the target variable remained similar in both samples.

Resampling and standardization

The data received was highly unbalanced - companies which defaulted accounted for only around 5,87% of the entire dataset. I used SMOTE to create new synthetic observations in a training sample based on existing ones and thus balancing the dataset.

I also standarized the data so that it can be correctly used by different models.

Variables selection

| | Var_01 | Var_02 | Var_04 | Var_05 | Var_06 | Var_07 | Var_08 | Var_09 | Var_10 | Var_11 | Var_12 | Var_13 | Var_14 | Var_15 | Var_16 | Var_18 | Var_19 | Var_21 | Var_22 | Var_23 | Var_24 | Var_25 | Var_27 | Var_28 | Var_29 | Var_30 | Var_31 | Var_32 | Var_33 | Var_34 | Var_35 | Var_36 | Var_37 | Var_39 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Var_01 | 1.00 | 0.95 | -0.14 | 0.01 | 0.95 | -0.06 | 0.84 | 0.87 | 0.83 | 0.95 | 0.89 | 0.02 | 0.89 | -0.21 | 1.00 | 0.95 | 0.95 | 0.97 | 0.79 | -0.07 | 0.81 | -0.13 | -0.13 | 0.13 | 0.93 | 0.93 | -0.13 | 0.73 | 0.02 | 0.83 | 0.05 | 0.73 | 0.02 | 0.71 |
| Var_02 | 0.95 | 1.00 | -0.15 | 0.02 | 1.00 | -0.01 | 0.85 | 0.91 | 0.86 | 1.00 | 0.94 | -0.00 | 0.91 | -0.14 | 0.95 | 1.00 | 1.00 | 0.94 | 0.79 | -0.04 | 0.81 | -0.14 | -0.14 | 0.17 | 0.89 | 0.89 | -0.20 | 0.80 | 0.05 | 0.86 | 0.02 | 0.80 | 0.05 | 0.63 |
| Var_04 | -0.14 | -0.15 | 1.00 | -0.17 | -0.15 | 0.20 | -0.13 | -0.14 | -0.08 | -0.15 | -0.11 | -0.23 | -0.16 | 0.14 | -0.14 | -0.15 | -0.15 | -0.20 | -0.10 | 0.27 | -0.11 | 0.94 | 0.21 | 0.08 | -0.17 | -0.17 | -0.09 | -0.19 | -0.11 | -0.08 | -0.49 | -0.19 | -0.11 | 0.07 |
| Var_05 | 0.01 | 0.02 | -0.17 | 1.00 | 0.02 | 0.11 | -0.02 | -0.00 | -0.06 | 0.02 | -0.04 | 0.71 | -0.02 | 0.07 | 0.01 | 0.02 | 0.02 | 0.03 | -0.05 | -0.05 | -0.05 | -0.17 | -0.21 | -0.18 | -0.02 | -0.02 | -0.17 | 0.19 | 0.43 | -0.06 | 0.46 | 0.19 | 0.43 | -0.06 |
| Var_06 | 0.95 | 1.00 | -0.15 | 0.02 | 1.00 | -0.02 | 0.86 | 0.91 | 0.88 | 1.00 | 0.94 | 0.00 | 0.91 | -0.14 | 0.95 | 1.00 | 1.00 | 0.95 | 0.80 | -0.05 | 0.82 | -0.14 | -0.14 | 0.16 | 0.90 | 0.90 | -0.19 | 0.80 | 0.05 | 0.86 | 0.03 | 0.80 | 0.04 | 0.64 |
| Var_07 | -0.08 | -0.01 | 0.20 | 0.11 | -0.02 | 1.00 | 0.05 | 0.04 | 0.01 | -0.01 | 0.00 | -0.07 | -0.04 | 0.64 | -0.08 | -0.01 | -0.01 | -0.08 | 0.04 | 0.63 | 0.04 | 0.26 | 0.21 | -0.03 | -0.13 | -0.13 | -0.42 | 0.04 | 0.21 | 0.01 | -0.12 | 0.04 | 0.21 | -0.07 |
| Var_08 | 0.84 | 0.85 | -0.13 | -0.02 | 0.86 | 0.05 | 1.00 | 0.96 | 0.76 | 0.85 | 0.84 | -0.04 | 0.90 | -0.11 | 0.84 | 0.85 | 0.85 | 0.82 | 0.92 | 0.02 | 0.94 | -0.13 | -0.03 | 0.21 | 0.84 | 0.84 | -0.14 | 0.68 | -0.02 | 0.76 | -0.01 | 0.68 | -0.03 | 0.62 |
| Var_09 | 0.87 | 0.91 | -0.14 | -0.00 | 0.91 | 0.04 | 0.96 | 1.00 | 0.79 | 0.91 | 0.88 | -0.04 | 0.94 | -0.14 | 0.87 | 0.91 | 0.91 | 0.86 | 0.88 | -0.01 | 0.91 | -0.14 | -0.06 | 0.20 | 0.88 | 0.88 | -0.16 | 0.73 | -0.01 | 0.79 | -0.00 | 0.73 | -0.01 | 0.62 |
| Var_10 | 0.83 | 0.86 | -0.08 | -0.06 | 0.88 | 0.01 | 0.76 | 0.79 | 1.00 | 0.86 | 0.92 | -0.11 | 0.78 | -0.14 | 0.83 | 0.86 | 0.86 | 0.79 | 0.75 | 0.02 | 0.76 | -0.07 | -0.11 | 0.22 | 0.77 | 0.77 | -0.19 | 0.62 | -0.01 | 1.00 | -0.14 | 0.62 | -0.01 | 0.65 |
| Var_11 | 0.95 | 1.00 | -0.15 | 0.02 | 1.00 | -0.01 | 0.85 | 0.91 | 0.86 | 1.00 | 0.95 | -0.00 | 0.91 | -0.14 | 0.95 | 1.00 | 1.00 | 0.94 | 0.79 | -0.04 | 0.81 | -0.14 | -0.14 | 0.17 | 0.89 | 0.89 | -0.20 | 0.80 | 0.05 | 0.86 | 0.02 | 0.80 | 0.05 | 0.63 |
| Var_12 | 0.89 | 0.94 | -0.11 | -0.04 | 0.94 | 0.00 | 0.84 | 0.88 | 0.92 | 0.95 | 1.00 | -0.10 | 0.88 | -0.12 | 0.89 | 0.94 | 0.95 | 0.86 | 0.81 | -0.00 | 0.82 | -0.10 | -0.12 | 0.19 | 0.84 | 0.84 | -0.21 | 0.71 | -0.00 | 0.92 | -0.11 | 0.71 | -0.00 | 0.66 |
| Var_13 | 0.02 | -0.00 | -0.23 | 0.71 | 0.00 | -0.07 | -0.04 | -0.04 | -0.11 | -0.00 | -0.10 | 1.00 | -0.03 | -0.07 | 0.02 | -0.00 | -0.00 | 0.05 | -0.07 | -0.14 | -0.07 | -0.21 | -0.17 | -0.40 | 0.01 | 0.01 | 0.05 | 0.08 | 0.22 | -0.11 | 0.61 | 0.08 | 0.23 | -0.07 |
| Var_14 | 0.89 | 0.91 | -0.16 | -0.02 | 0.91 | -0.04 | 0.90 | 0.94 | 0.78 | 0.91 | 0.88 | -0.03 | 1.00 | -0.10 | 0.89 | 0.91 | 0.91 | 0.88 | 0.83 | -0.07 | 0.85 | -0.16 | -0.09 | 0.27 | 0.91 | 0.91 | -0.14 | 0.72 | -0.02 | 0.78 | 0.02 | 0.72 | -0.02 | 0.62 |
| Var_15 | -0.21 | -0.14 | 0.14 | 0.07 | -0.14 | 0.64 | -0.11 | -0.14 | -0.14 | -0.14 | -0.12 | -0.07 | -0.10 | 1.00 | -0.21 | -0.14 | -0.14 | -0.20 | -0.10 | 0.42 | -0.11 | 0.22 | 0.13 | 0.07 | -0.25 | -0.25 | -0.41 | -0.08 | 0.11 | -0.14 | -0.09 | -0.06 | 0.11 | -0.18 |
| Var_16 | 1.00 | 0.95 | -0.14 | 0.01 | 0.95 | -0.08 | 0.84 | 0.87 | 0.83 | 0.95 | 0.89 | 0.02 | 0.89 | -0.21 | 1.00 | 0.95 | 0.95 | 0.97 | 0.79 | -0.07 | 0.81 | -0.13 | -0.13 | 0.13 | 0.93 | 0.93 | -0.13 | 0.73 | 0.02 | 0.83 | 0.05 | 0.73 | 0.02 | 0.71 |
| Var_18 | 0.95 | 1.00 | -0.15 | 0.02 | 1.00 | -0.01 | 0.85 | 0.91 | 0.86 | 1.00 | 0.94 | -0.00 | 0.91 | -0.14 | 0.95 | 1.00 | 1.00 | 0.94 | 0.79 | -0.04 | 0.81 | -0.14 | -0.14 | 0.17 | 0.89 | 0.89 | -0.20 | 0.80 | 0.05 | 0.86 | 0.02 | 0.80 | 0.05 | 0.63 |
| Var_19 | 0.95 | 1.00 | -0.15 | 0.02 | 1.00 | -0.01 | 0.85 | 0.91 | 0.86 | 1.00 | 0.95 | -0.00 | 0.91 | -0.14 | 0.95 | 1.00 | 1.00 | 0.94 | 0.79 | -0.04 | 0.81 | -0.14 | -0.14 | 0.17 | 0.89 | 0.89 | -0.20 | 0.80 | 0.05 | 0.86 | 0.05 | 0.80 | 0.05 | 0.63 |
| Var_21 | 0.97 | 0.94 | -0.20 | 0.03 | 0.95 | -0.08 | 0.82 | 0.86 | 0.79 | 0.94 | 0.86 | 0.05 | 0.88 | -0.20 | 0.97 | 0.94 | 0.94 | 1.00 | 0.75 | -0.09 | 0.78 | -0.19 | -0.14 | 0.11 | 0.93 | 0.93 | -0.13 | 0.76 | 0.04 | 0.79 | 0.09 | 0.76 | 0.04 | 0.59 |
| Var_22 | 0.79 | 0.79 | -0.10 | -0.05 | 0.80 | 0.04 | 0.92 | 0.88 | 0.75 | 0.79 | 0.81 | -0.07 | 0.83 | -0.10 | 0.79 | 0.79 | 0.79 | 0.75 | 1.00 | 0.11 | 0.98 | -0.10 | 0.03 | 0.25 | 0.78 | 0.78 | -0.13 | 0.59 | -0.05 | 0.75 | -0.07 | 0.59 | -0.06 | 0.62 |
| Var_23 | -0.07 | -0.04 | 0.27 | -0.05 | -0.05 | 0.63 | 0.02 | -0.01 | 0.02 | -0.04 | -0.00 | -0.14 | -0.07 | 0.42 | -0.07 | -0.04 | -0.04 | -0.09 | 0.11 | 1.00 | 0.09 | 0.31 | 0.47 | 0.01 | -0.12 | -0.12 | -0.22 | -0.05 | 0.02 | 0.02 | -0.31 | -0.05 | 0.01 | -0.03 |
| Var_24 | 0.81 | 0.81 | -0.11 | -0.05 | 0.82 | 0.04 | 0.94 | 0.91 | 0.76 | 0.81 | 0.82 | -0.07 | 0.85 | -0.11 | 0.81 | 0.81 | 0.81 | 0.78 | 0.98 | 0.09 | 1.00 | -0.10 | 0.01 | 0.24 | 0.80 | 0.80 | -0.13 | 0.61 | -0.05 | 0.76 | -0.06 | 0.61 | -0.05 | 0.63 |
| Var_25 | -0.13 | -0.14 | 0.94 | -0.17 | -0.14 | 0.26 | -0.13 | -0.14 | -0.07 | -0.14 | -0.10 | -0.21 | -0.16 | 0.22 | -0.13 | -0.14 | -0.14 | -0.19 | -0.10 | 0.31 | -0.10 | 1.00 | 0.20 | 0.02 | -0.18 | -0.18 | -0.13 | -0.18 | -0.10 | -0.07 | -0.47 | -0.18 | -0.10 | 0.06 |
| Var_27 | -0.13 | -0.14 | 0.21 | -0.21 | -0.14 | 0.21 | -0.03 | -0.06 | -0.11 | -0.14 | -0.12 | -0.17 | -0.09 | 0.13 | -0.13 | -0.14 | -0.14 | -0.14 | 0.03 | 0.47 | 0.01 | 0.20 | 1.00 | 0.05 | -0.11 | -0.11 | 0.36 | -0.17 | -0.25 | -0.11 | -0.28 | -0.17 | -0.26 | -0.07 |
| Var_28 | 0.13 | 0.17 | 0.08 | -0.18 | 0.16 | 0.21 | 0.20 | 0.22 | 0.17 | 0.17 | 0.19 | -0.40 | 0.27 | 0.07 | 0.13 | 0.17 | 0.17 | 0.11 | 0.25 | 0.01 | 0.24 | 0.02 | 0.05 | 1.00 | 0.17 | 0.17 | -0.08 | 0.12 | -0.05 | 0.22 | -0.18 | 0.12 | -0.05 | 0.17 |
| Var_29 | 0.93 | 0.89 | -0.17 | -0.02 | 0.90 | -0.13 | 0.84 | 0.88 | 0.77 | 0.89 | 0.84 | 0.01 | 0.91 | -0.25 | 0.93 | 0.89 | 0.89 | 0.93 | 0.78 | -0.12 | 0.80 | -0.18 | -0.11 | 0.17 | 1.00 | 1.00 | -0.02 | 0.69 | -0.02 | 0.77 | 0.05 | 0.69 | -0.02 | 0.64 |
| Var_30 | 0.93 | 0.89 | -0.17 | -0.02 | 0.90 | -0.13 | 0.84 | 0.88 | 0.77 | 0.89 | 0.84 | 0.01 | 0.91 | -0.25 | 0.93 | 0.89 | 0.89 | 0.93 | 0.78 | -0.12 | 0.80 | -0.18 | -0.11 | 0.17 | 1.00 | 1.00 | -0.02 | 0.69 | -0.02 | 0.77 | 0.05 | 0.69 | -0.02 | 0.64 |
| Var_31 | -0.13 | -0.20 | -0.09 | -0.17 | -0.19 | -0.42 | -0.14 | -0.16 | -0.19 | -0.20 | -0.21 | 0.05 | -0.14 | -0.41 | -0.13 | -0.20 | -0.20 | -0.13 | -0.13 | -0.22 | -0.13 | -0.13 | 0.36 | -0.08 | -0.02 | -0.02 | 1.00 | -0.22 | -0.25 | -0.19 | 0.09 | -0.22 | -0.27 | -0.10 |
| Var_32 | 0.73 | 0.80 | -0.19 | 0.19 | 0.80 | 0.04 | 0.68 | 0.73 | 0.62 | 0.80 | 0.71 | 0.08 | 0.72 | -0.08 | 0.73 | 0.80 | 0.80 | 0.76 | 0.59 | -0.05 | 0.61 | -0.18 | -0.17 | 0.12 | 0.69 | 0.69 | -0.22 | 1.00 | 0.21 | 0.62 | 0.17 | 1.00 | 0.21 | 0.45 |
| Var_33 | 0.02 | 0.05 | -0.11 | 0.43 | 0.05 | 0.21 | -0.02 | -0.01 | -0.01 | 0.05 | -0.00 | 0.22 | -0.02 | 0.11 | 0.02 | 0.05 | 0.05 | 0.04 | -0.05 | 0.02 | -0.05 | -0.10 | -0.25 | -0.05 | -0.02 | -0.02 | -0.26 | 0.21 | 1.00 | -0.01 | 0.32 | 0.21 | 0.99 | -0.05 |
| Var_34 | 0.83 | 0.86 | -0.08 | -0.06 | 0.86 | 0.01 | 0.76 | 0.79 | 1.00 | 0.86 | 0.92 | -0.11 | 0.78 | -0.14 | 0.83 | 0.86 | 0.86 | 0.79 | 0.75 | 0.02 | 0.76 | -0.07 | -0.11 | 0.22 | 0.77 | 0.77 | -0.19 | 0.62 | -0.01 | 1.00 | -0.14 | 0.62 | -0.01 | 0.65 |
| Var_35 | 0.05 | 0.02 | -0.49 | 0.46 | 0.03 | -0.12 | -0.01 | -0.00 | -0.14 | 0.02 | -0.11 | 0.61 | 0.02 | -0.09 | 0.05 | 0.02 | 0.02 | 0.09 | -0.07 | -0.31 | -0.06 | -0.47 | -0.28 | -0.18 | 0.05 | 0.05 | 0.09 | 0.17 | 0.32 | -0.14 | 1.00 | 0.17 | 0.32 | -0.10 |
| Var_36 | 0.73 | 0.80 | -0.19 | 0.19 | 0.80 | 0.04 | 0.68 | 0.73 | 0.62 | 0.80 | 0.71 | 0.08 | 0.72 | -0.08 | 0.73 | 0.80 | 0.80 | 0.76 | 0.59 | -0.05 | 0.61 | -0.18 | -0.17 | 0.12 | 0.69 | 0.69 | -0.22 | 1.00 | 0.21 | 0.62 | 0.17 | 1.00 | 0.21 | 0.45 |
| Var_37 | 0.02 | 0.05 | -0.11 | 0.43 | 0.04 | 0.21 | -0.03 | -0.01 | -0.01 | 0.05 | -0.00 | 0.23 | -0.02 | 0.11 | 0.02 | 0.05 | 0.05 | 0.04 | -0.06 | 0.01 | -0.05 | -0.10 | -0.26 | -0.05 | -0.02 | -0.02 | -0.27 | 0.21 | 0.99 | -0.01 | 0.32 | 0.21 | 1.00 | -0.05 |
| Var_39 | 0.71 | 0.63 | 0.07 | -0.06 | 0.64 | -0.07 | 0.62 | 0.62 | 0.65 | 0.63 | 0.66 | -0.07 | 0.62 | -0.18 | 0.71 | 0.63 | 0.63 | 0.59 | 0.62 | -0.03 | 0.63 | 0.06 | -0.07 | 0.17 | 0.64 | 0.64 | -0.10 | 0.45 | -0.05 | 0.65 | -0.10 | 0.45 | -0.05 | 1.00 |

Correlation plot for X_train after resampling, outliers handling and standardization

Based on the Pearson correlation plot, I see that some of the variables are highly correlated. Due to time and resources constraints, I decided to drop selected variables manually (instead of, i.e. performing PCA) based on their correlation with other variables and economic interpretation. All variables deleted below had very high (coefficient of almost 1) correlation with kept variables and indicated similar positions in the balance sheet such as IFRS Assets Total vs. Assets Total or Tangible Net Worth and Eff Tang Net Worth Actual.

Deleted variables:
- Var_06, Var_11, Var_18, Var_19 - keeping Var_02
- Var_32 - keeping Var_36
- Var_10 - keeping Var_34
- Var_16 - keeping Var_01

Additionally, I added 6 dummy variables indicating the year (2016 to 2021 and 2015 as default) in which each observation was observed.

Finally, the training dataset consisted of 196252 observations with 40 variables including 34 numeric variables on financial data and 6 dummy year variables. The dataset was balanced in terms of target variable thanks to SMOTE.

## Modelling

For all models, I used F1 score as the primary evaluation metric for tuning the parameters. F1 is the harmonic mean of precision and recall and I used it because it is more suitable for unbalanced data than simple accuracy or ROC AUC score (I rebalanced training data, but will still eventually evaluate performance on unbalanced test data).
To properly evaluate models' results and avoid overfitting, each time I used five fold cross-validation (repeated three times in selected algorithms).

Logistic regression
Logistic regression (logit) is one of the most basic models for classification and was used here as a benchmark for more sophisticated models. I achieved an F1 score of 0.7963 which is okay, but not very good.

Ridge regression
Ridge regression is a method of regularization which imposes restrictions on less important parameters that bring them close to 0. I tuned in penalization parameter alpha and the best F1 score of 0.7936 was obtained for alpha = 1. Therefore I were not able to improve results compared to logit.

K nearest neighbors (KNN)
KNN is a non-parametric model which uses proximity to make classifications about the grouping of an individual data point. I tuned in the number of neighbors parameter (how many neighbors of an observation will be checked to determine the classification). I obtained the best result of F1 score of 0.9698 for 11 neighbors. It is a significant improvement compared to logit and ridge and overall a very good result. However, I should be careful and check it later on testing data as it might be due to overfitting.

XGBoost
XGboost is a type of gradient boosting i.e. a procedure, where final prediction is created by combining results of many decision trees which are dependent on the previous ones. I tuned the learning rate, number of estimators, subsample and maximum depth parameters. I obtained the best F1 score of 0.9677 for learning_rate = 0.1, max_depth = 9, n_estimators = 500 and subsample = 0.8. Thus, the results were very similar to KNN.

# Comparison and conclusions

Let's recap the performance of models on training data and check and compare their performance on testing data.

| Sample & performance measure | Logit | Ridge | KNN | XGBoost |
|---|---|---|---|---|
| Training F1 | 0.7963 | 0.7936 | 0.9698 | **0.9677** |
| Testing F1 | 0.2392 | 0.2365 | 0.2524 | **0.4483** |
| Testing ROC AUC | 0.6730 | 0.6594 | 0.7205 | **0.7168** |

On training data, the best models in terms of F1 score were KNN and XBoost with F1 score of around 0.97 for both. Logit and Ridge performed much worse and both had an F1 score of around 0.8.

On testing data, unfortunately the results in terms of F1 were very bad, most probably as a result of some error - maybe in the prediction procedure itself (which is rather not probable but what I hope for), but also possibly for example in resampling or optimization of parameters as the ROC AUC score were poor, but not terrible which could indicate that the negative class was predicted okay.

If I must choose a model for prediction, I would pick XGBoost, as it has significantly higher performance on the testing sample than other methods but I acknowledge that it is still poor and should be investigated further with checking classification matrix etc.

Further analysis

I acknowledge that due to lack of time and processing power the analysis contains some shortcomings. Given more time and resources, the author would implement the following:
- First and foremost, investigating the origin of such terrible results on testing data.
- Better handling of missing values and outliers, for example investigating the cases for each variable separately
- Performing PCA on all or selected highly correlated variables and treating them as predictors to reduce dimension of the dataset. I did not do that as this would require a significant amount of additional time put into ensuring economic interpretation of obtained principal components.
- More robust validation, such as increasing the number of folds in cross-validation and optimizing parameters based on multiple evaluation metrics, not just F1.
- Implementing support vector machine (SVM) and neural networks models which require more time spent on computation and parameters tuning.