# Configuration Manual

MSc Research Project
FinTech

## Saumya Chaudhary

Student ID: x21175365

School of Computing
National College of Ireland

Supervisor: Brian Byrne

| | |
|---|---|
| **Student Name:** | Saumya Chaudhary |
| **Student ID:** | x21175365 |
| **Programme:** | FinTech |
| **Year:** | 2023 |
| **Module:** | MSc Research Project |
| **Supervisor:** | Brian Byrne |
| **Submission Due Date:** | 14/08/2023 |
| **Project Title:** | Configuration Manual |
| **Word Count:** | 420 |
| **Page Count:** | 5 |

| **Signature:** | Saumya Chaudhary |
|---|---|
| **Date:** | 17th September 2023 |

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:**

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies). | ☐ |
| **Attach a Moodle submission receipt of the online project submission**, to each project (including multiple copies). | ☐ |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. | ☐ |

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

| **Office Use Only** | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# Configuration Manual

## Saumya Chaudhary
x21175365

# 1 Environment Setup

The studies were conducted in a Google Colab environment (Bisong and Bisong (2019)) utilizing a Jupyter Notebook (Kluyver et al. (2016)) and the Python programming language(Sanner et al. (1999)). Table 2 shows the various tools used with their purpose.

| Sl. No. | Tools and Technologies | Purpose |
|---------|------------------------|---------|
| 1 | Microsoft Excel | *For the purpose of storing and retrieving Taiwanese bankruptcy data in CSV format* |
| 2 | Google Colaboratory | *Free, cloud-based Jupyter notebook for collaborative Python programming and experimentation* |
| 3 | Jupyter Notebook | *To implement Python Codes* |
| 4 | Python | *For carrying out a high-level coding task for machine learning* |

Table 1: Various Tools and Technologies, and Purpose

Following Python Packages were used

| Name | Description |
|------|-------------|
| Pandas | *Data processing* |
| Numpy | *Linear algebra with features* |
| Seaborn, Pyplot and Matplotlib | *Data Visualization* |
| Sci-kit learn | *For Data pre-processing and model implementation* |
| Imbalanced learn | *For Data resampling and ensemble model implementation* |

Table 2: Various Tools and Technologies, and Purpose

# 2    Data Analysis

## 2.1    Data Upload

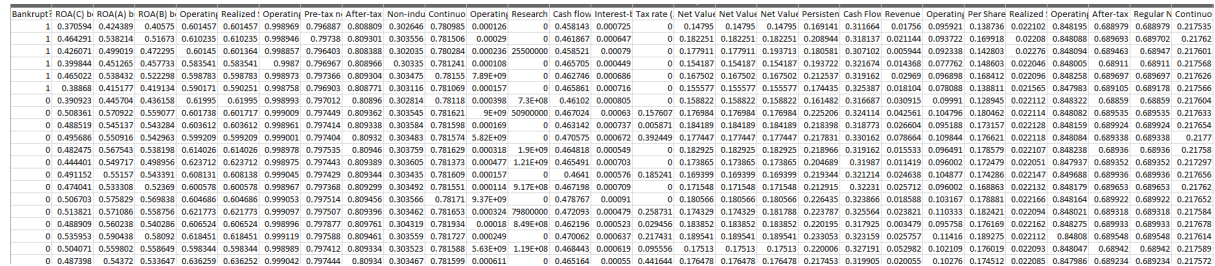The data was collected in CSV format as shown below



Figure 1: Data in Excel

## 2.2    Checking for null values and duplicates

It's crucial to check if there are any missing values in the dataset as this can degrade the model quality. Also, checking for duplicates is an important phase. The steps include

Step 1: df.info() gives the null value information Step 2: df.duplicated() gives number of duplicates values.

## 2.3    Feature Selection

The features were chosen using Pearson's correlation co-efficients as shown below. This create a set of features by dropping highly correlated independent variables. Based on the literature, it is found that $> |0.7|$ is assumed to be strong correlation (Dormann et al. (2013)).

## 2.4    Data Normalization

All the variables were scaled using min-max scaler and standard scaler as shown below.

## 2.5    Dimensionality Reduction

Principal Component Analysis (PCA) is a widely used technique for reducing the dimensionality of high-dimensional data while preserving as much variance as possible (Abdi and Williams (2010)). It entails converting the initial variables into a new collection of uncorrelated variables (principal components) that account for the majority of the variability in the data. In this study, the top 8 principal components accounted for 99% of the variation in the dataset.

## 2.6    Train and Test Split

Splitting data into training and test set are important as they will be used to train and evaluate the models. The 70% of the data were used to train the model while 30% to test the model using stratified sampling strategy.

| | Bankrupt? | ROA(C) before interest and depreciation before interest | ROA(A) before interest and % after tax | ROA(B) before interest and depreciation after tax | Operating Gross Margin | Realized Sales Gross Margin | Operating Profit Rate | Pre-tax net Interest Rate |
|---|---|---|---|---|---|---|---|---|
| Net Income to Total Assets | -0.2343 | 0.9244 | 0.9916 | 0.9326 | 0.5054 | 0.5034 | 0.7979 | 0.9066 |
| Total assets to GNP price | 0.0551 | -0.0068 | 0.0104 | 0.0030 | -0.2580 | -0.2581 | 0.0156 | 0.0624 |
| No-credit Interval | -0.1117 | 0.1893 | 0.2126 | 0.1807 | 0.0974 | 0.0965 | 0.1567 | 0.1925 |
| Gross Profit to Sales | -0.1476 | 0.5062 | 0.4965 | 0.5157 | 1.0000 | 0.9991 | 0.6856 | 0.5870 |
| Net Income to Stockholder's Equity | -0.2191 | 0.8989 | 0.9694 | 0.9060 | 0.4505 | 0.4483 | 0.7655 | 0.8575 |
| Liability to Equity | 0.2007 | -0.2409 | -0.2357 | -0.2456 | -0.3816 | -0.3822 | -0.2377 | -0.3301 |
| Degree of Financial Leverage (DFL) | -0.1471 | 0.2199 | 0.2355 | 0.2282 | 0.0416 | 0.0415 | 0.2304 | 0.2430 |
| Interest Coverage Ratio (Interest expense to EBIT) | -0.1275 | 0.1836 | 0.1910 | 0.1930 | 0.0159 | 0.0151 | 0.1979 | 0.1913 |
| Equity to Liability | -0.2172 | 0.2475 | 0.2425 | 0.2523 | 0.3858 | 0.3864 | 0.2434 | 0.3365 |

Figure 2: Correlation Between a Few Variables



Transform data using min-max scaler and standard scaler

As we can see, there are many outliers in the dataset, we can perform min-max scaling.

A `Min-Max Scaling` is typically done via the foloowing equation:

$$X_{norm} = \frac{X_i - X_{min}}{X_{max} - X_{min}}$$

$X_i$ is the $i^{th}$ sample of dataset.

We will also do experiment with standardized or Z-score normalized data. The features will be rescaled as a result of standardization (or Z-Score normalization), and they will then possess the characteristics of a conventional normal distribution with:

$$\mu = 0$$

And

$$\sigma = 1$$

Where $\mu$ is the mean(average) and $\sigma$ is the standard deviation from the mean; standard scores (also called **Z** scores) of the sampels are calculated as follows:

$$z = \frac{x - \mu}{\sigma}$$

Reference: https://colab.research.google.com/github/JL1829/turbo-funicular/blob/master/_notebooks/2020-09-19-About%20Feature%20Scaling%20and%20Normalization.ipynb#scrollTo=TetP0B43NPye

```
[ ]  min_max = MinMaxScaler()
     std_scl = StandardScaler()

     X_min_max = min_max.fit_transform(X_non_corr)
     X_std_scl = std_scl.fit_transform(X_non_corr)
```
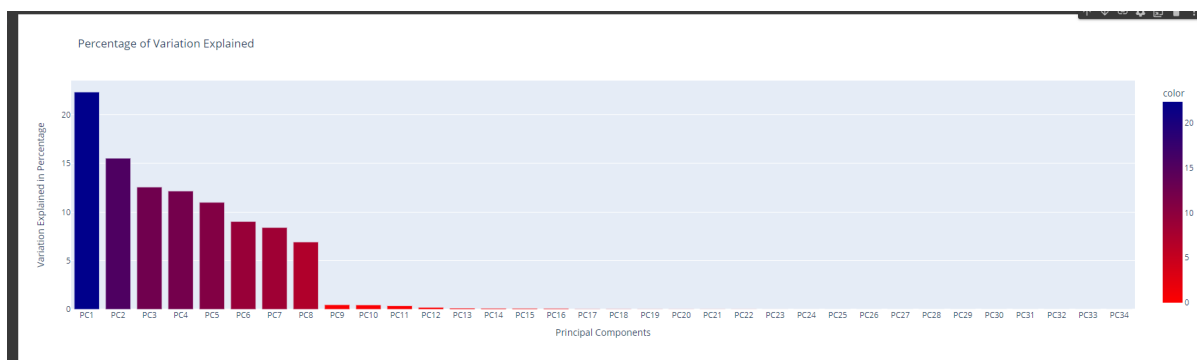
Figure 3: Scaling

Figure 4: PCA

## 2.7 Resampling

As the data is biased towards non-bankrupt events, the resampling strategies including Synthetic Minority Oversampling Technique (SMOTE) (Chawla et al. (2002); Smiti and Soui (2020)), Random Oversampling (ROS) (Mohammed et al. (2020)), and Random Undersampling (RUS) (Wang and Liu (2021); Liu et al. (2008)) have been implemented. The SMOTE and Over-sampling created new data points for minority class, i.e., class with bankruptcy events.

# References

Abdi, H. and Williams, L. J. (2010). Principal component analysis, *Wiley interdisciplinary reviews: computational statistics* **2**(4): 433–459.

Bisong, E. and Bisong, E. (2019). Google colaboratory, *Building machine learning and deep learning models on google cloud platform: a comprehensive guide for beginners* pp. 59–64.

Chawla, N. V., Bowyer, K. W., Hall, L. O. and Kegelmeyer, W. P. (2002). Smote: synthetic minority over-sampling technique, *Journal of artificial intelligence research* **16**: 321–357.

Dormann, C. F., Elith, J., Bacher, S., Buchmann, C., Carl, G., Carré, G., Marquéz, J. R. G., Gruber, B., Lafourcade, B., Leitão, P. J. et al. (2013). Collinearity: a review of methods to deal with it and a simulation study evaluating their performance, *Ecography* **36**(1): 27–46.

Kluyver, T., Ragan-Kelley, B., Pérez, F., Granger, B. E., Bussonnier, M., Frederic, J., Kelley, K., Hamrick, J. B., Grout, J., Corlay, S. et al. (2016). Jupyter notebooks-a publishing format for reproducible computational workflows., *Elpub* **2016**: 87–90.

Liu, X.-Y., Wu, J. and Zhou, Z.-H. (2008). Exploratory undersampling for class-imbalance learning, *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* **39**(2): 539–550.

Mohammed, R., Rawashdeh, J. and Abdullah, M. (2020). Machine learning with over-sampling and undersampling techniques: overview study and experimental results, *2020*

*11th international conference on information and communication systems (ICICS)*, IEEE, pp. 243–248.

Sanner, M. F. et al. (1999). Python: a programming language for software integration and development, *J Mol Graph Model* **17**(1): 57–61.

Smiti, S. and Soui, M. (2020). Bankruptcy prediction using deep learning approach based on borderline smote, *Information Systems Frontiers* **22**: 1067–1083.

Wang, H. and Liu, X. (2021). Undersampling bankruptcy prediction: Taiwan bankruptcy data, *Plos one* **16**(7): e0254030.