Bach Quang Minh

Le Ngoc Thach

Nguyen Nhat Quang

Jon Von Neumann Institute – ICT2017

## HACKATHON II - 2018

## Problem statement

Measure the default level of a company based on its financial profile.

## Data exploration

The whole data set contains 12 features and 2 classes: 'default'-1 and 'not-default'-0. We identify that this is a binary classification problem.

Here is the distribution of the data class:

```
In [275]: df.default.value_counts()
Out[275]:
0.0    357
1.0     20
Name: default, dtype: int64
```

As seen, the data suffered from imbalanced-class issue. The class percentage of 1 and 0 are 5.3% and 94.7% respectively.

We also notice that there is one company was having an unusual 'Account & Notes receivables' value:

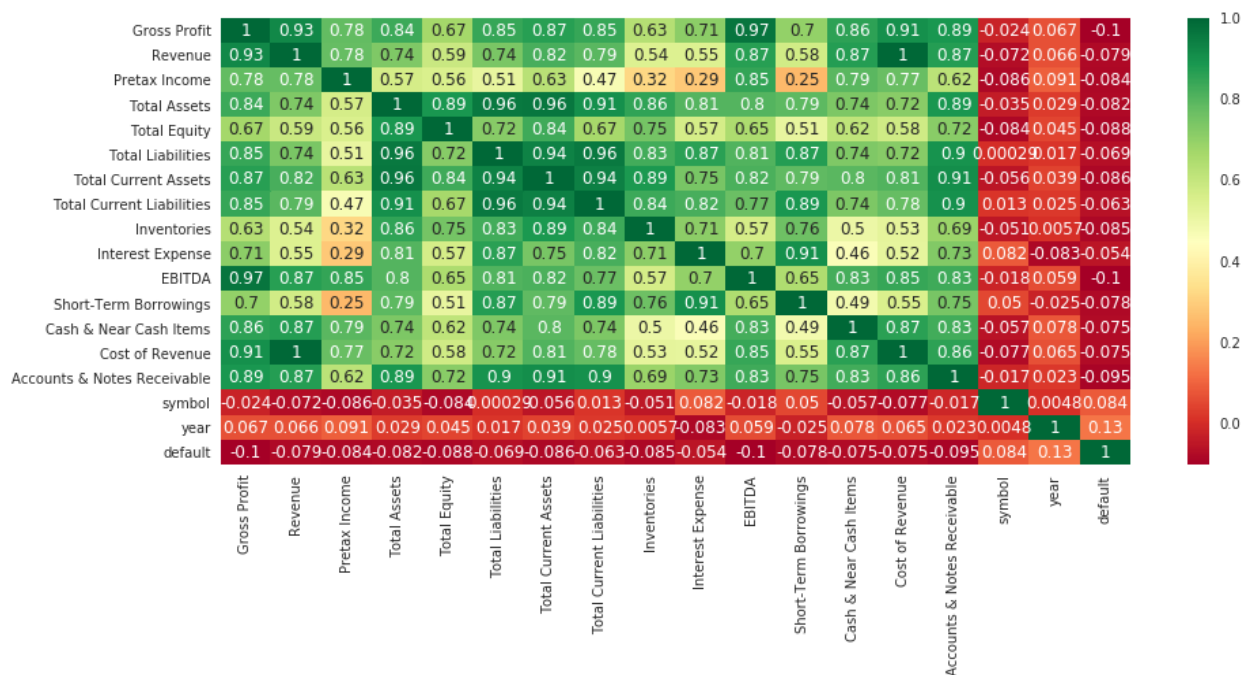| Total Current Asse | Total Current Liabiliti | Inventorie | Interest Expen | EBITDA | Short-Term Borrowing | Cash & Near Cash Iten | Cost of Revent | Accounts & Notes Receivabl | symb | ye | defa |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1038865.031 | 1008993.567 | 562910.7612 | 49837.84126 | 92549.30463 | 392909.4739 | 17965.52508 | 761840.5745 | 310805.2571 | 41 | 2013 | 0 |
| 295739.5764 | 236844.3272 | 63834.89325 | 4894.87626 | 36754.56997 | 70334.00531 | 112384.1551 | 437038.4766 | 101629.0227 | 89 | 2014 | 0 |
| 684080.2622 | 449868.4911 | 383294.9369 | 12852.54206 | 19186.0783 | 232444.5774 | 17215.35261 | 182182.1609 | 192038.6579 | 8 | 2012 | 0 |
| 665498.9886 | 344484.3875 | 407121.3734 | 2478.295 | -7111.191917 | 101620 | 587.979562 | 13771.84937 | 91 | 63 | 2014 | 0 |
| 607741.414 | 459906.2602 | 47537.99559 | 2296.149575 | 61499.48046 | 34410.81486 | 204385.8563 | 573892.4953 | 156344.8002 | 92 | 2015 | 0 |

Since, this is the company data, we decided to drop this row before further analysis.

## Proposed method

Since most columns from the data are real, continuous values. We decided to use Logistic Regression – a classification algorithm which could deal with continuous data.

The measure that we use for evaluation is the True Positive rate.

**Feature selection**



After studying carefully the correlation matrix and combining with our understanding about these finance indicators, we decided to omit the following features:

1. Total equity (because this is just an index calculated from two others indices)
2. EBITDA
3. Cost of Revenue (= revenue – gross profit)
4. Symbol (Just an encoded for each company)
5. Year (How well a company perform in a year is reflected by other index)

We will keep the following features for our first attempt:

'Gross Profit'

'Revenue'

'Pretax Income'

'Total Assets'

'Total Liabilities'

'Total Current Assets'

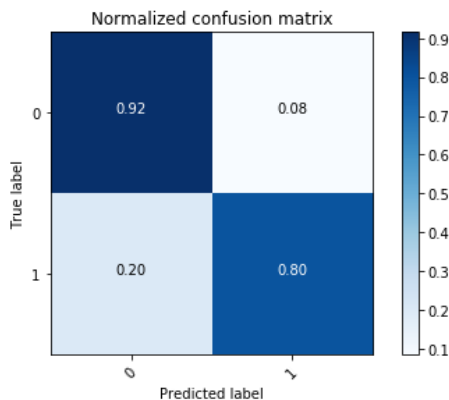'Total Current Liabilities','Inventories'

'Interest Expense'

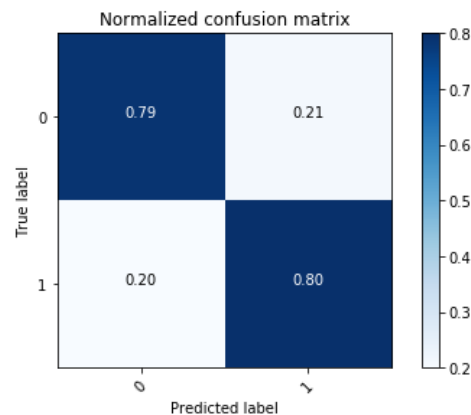'Short-Term Borrowings'

'Cash & Near Cash Items'

We use the train set and split it into train set and validation set with the ratio of 5:1

Here is the result:

Confusion matrix on the train set | Confusion matrix on the test set

Normalized confusion matrix



To know if our features are statistically significant, we investigate these statistics:

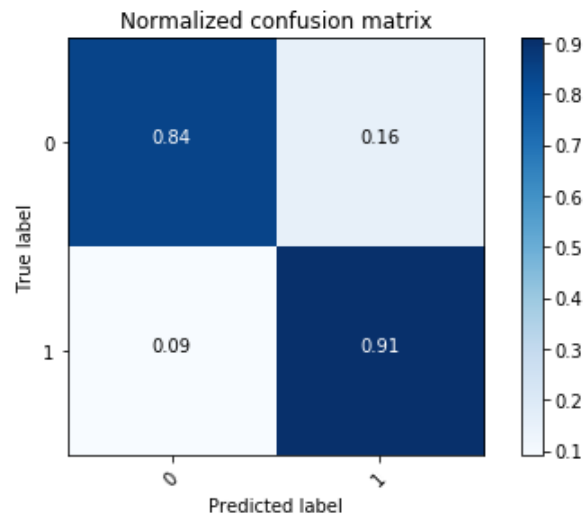| | Coefficients | Standard Errors | t values | Probabilites |
|---|---|---|---|---|
| 0 | 7.378570e-08 | 4.062806e-07 | 0.181613 | 0.856010 |
| 1 | -2.379103e-05 | 4.550089e-08 | -522.869570 | 0.000000 |
| 2 | 3.588993e-06 | 3.056883e-07 | 11.740693 | 0.000000 |
| 3 | -9.496157e-05 | 4.405679e-08 | -2155.435292 | 0.000000 |
| 4 | 1.640583e-06 | 6.699473e-08 | 24.488238 | 0.000000 |
| 5 | -5.962702e-07 | 1.288151e-07 | -4.628884 | 0.000005 |
| 6 | -8.132438e-06 | 9.960979e-08 | -81.642959 | 0.000000 |
| 7 | 9.204492e-06 | 1.066420e-07 | 86.312097 | 0.000000 |
| 8 | -2.870650e-06 | 9.008965e-07 | -3.186437 | 0.001593 |
| 9 | -6.825084e-05 | 1.674233e-07 | -407.654333 | 0.000000 |
| 10 | -4.264054e-06 | 2.554688e-07 | -16.691095 | 0.000000 |
| 11 | 6.646292e-06 | 2.097341e-02 | 0.000317 | 0.999747 |

The result shows that some features does not have small p-values. We decided to seek for further feature selections.

First, we use the train-validation ratio as 0.3 and construct **a for loop** for **the number of features** to select the model with **the highest true positive values.**

The loop gives us these features:

['Gross Profit', 'Revenue', 'Pretax Income', 'Total Assets', 'Total Current Assets']

This is the confusion matrix for the validation set

Normalized confusion matrix

|  | Predicted 0 | Predicted 1 |
|---|---|---|
| True 0 | 0.84 | 0.16 |
| True 1 | 0.09 | 0.91 |

The result obtained is higher and the features are suitable to our understanding.

We decided to use these features to train the model and produce the result for the test set.