# Predicting Corporate Bankruptcy Using Financial Ratios via Machine Learning

Nathan Silverglate, Brandon Brooks, Mayuresh Deolekar

Data Science for Business - Technical
Fall 2021

**NYU | STERN**

# Business Case

**Problem Statement**

- Business customers currently don't have a tool that allows them to assess whether or not a publicly-traded company is at risk of bankruptcy before buying company stock based on commonly available business financial ratios

# Business Case

## Our Solution

- A corporate bankruptcy prediction tool that gives customers the insights to properly assess a company's financial health before and while investing their capital

# Business Case

**Our predictive tool will help customers:**

- Collect the necessary data to make better and more informed decisions
- Gain critical insight to at-risk companies in the same industry and notify them of financial risk across that industry
- Highlight critical business issues for at-risk companies that the management team might be withholding
- Maximize their ROI
- Further hold public companies accountable

NYU | STERN

# Business Case

**Data mining is extremely useful in discovering patterns and correlations for various players in the finance industry. Here are few alternative use cases for our tool**:

- Financial institutions
- Professional investors
- Governments
- Academic researchers

# Data

- Data source from Kaggle
- Data collected from *Taiwan Economic Journal* for years 1999-2009
- Bankruptcy is defined by regulations of Taiwanese stock exchange

# Data

- Initial inspection yielded data frame of 96 columns with 6,819 entries
- 93 of columns are continuous numerical features corresponding to values including:
    - Operating Gross Margin
    - Realized Sales Gross Margin
    - Operating Profit Rate
    - Pre-tax net Interest Rate
    - After-tax net Interest Rate

NYU | STERN

# Data

- Two of the columns are binary features:
  - Liability-Assets Flag
  - Net Income Flag
  - Both of these features contained same values for all entries and were later dropped
- Final column is a binary classifier column corresponding to Bankruptcy Occurred?
  - Yes = 1
  - No = 0

# Data

- Dataset lacks any temporal features, prevents us from knowing details regarding prediction timing

- Dataset also lacks unique identifiers which prevents us from knowing if multiple instances represent the same company

- The only fixed value columns are the two binary flag features that are to be removed

- All other features are numerical and are actual business financial ratio metrics

# Data

- Upon initial inspection, data contained zero NaN values
- Since there are 95 features in dataset, needs to simplify to a more usable number of features
- Performed **Model Based Feature Selection** using **Random Forest Classifier** to judge importance of each feature and how many to use based on the **sensitivity score** of confusion matrix
- Found that **15 features** was optimal

```
Sensitivity Score: 0.164    number of features: 5
Sensitivity Score: 0.230    number of features: 10
Sensitivity Score: 0.262    number of features: 15
Sensitivity Score: 0.262    number of features: 20
Sensitivity Score: 0.213    number of features: 25
Sensitivity Score: 0.230    number of features: 30
Sensitivity Score: 0.230    number of features: 35
Sensitivity Score: 0.246    number of features: 40
```

# Data

The 15 selected features are as follows:

- Operating Expense Rate
- Interest-Bearing Debt Interest Rate
- Net Value Per Share (B)
- Persistent EPS in the Last Four Seasons
- Net Value Growth Rate
- Total Asset Return Growth Rate
- Interest Expense Ratio
- Borrowing Dependency

- Net Profit Before Tax/Paid-in Capital
- Cash/Total Assets
- Inventory/Working Capital
- Working Capital/Equity
- Total Assets to GNP Price
- Net Income to Stockholder's Equity
- Degree of Financial Leverage (DFL)

NYU | STERN

# Data

- Next investigated the distribution of target column values

- Distribution was quite skewed
  - Bankrupt = No: 6,599
  - Bankrupt = Yes: 220

- Ratio: 3.23% of entries went bankrupt



NYU | STERN

# Data

## Feature Distribution

- Data is both heavy on outliers or contains most values in a single "bin"

# Data

## Total Asset Return Growth Rate

- Has normal distribution but

  influenced by outliers



```
Total Asset Return Growth Rate Ratio
(0.2, 0.3]                              6815
(0.3, 0.4]                                 2
(0.9, 1.0]                                 1
(-0.001, 0.1]                              1
(0.8, 0.9]                                 0
(0.7, 0.8]                                 0
(0.6, 0.7]                                 0
(0.5, 0.6]                                 0
(0.4, 0.5]                                 0
(0.1, 0.2]                                 0
dtype: int64
```

```
Rows with outliers: 6819
Rows withou outliers: 6767
information lost = 52 rows
<matplotlib.axes._subplots.AxesSubplot at 0x7f9e1b736e10>
```

# Data

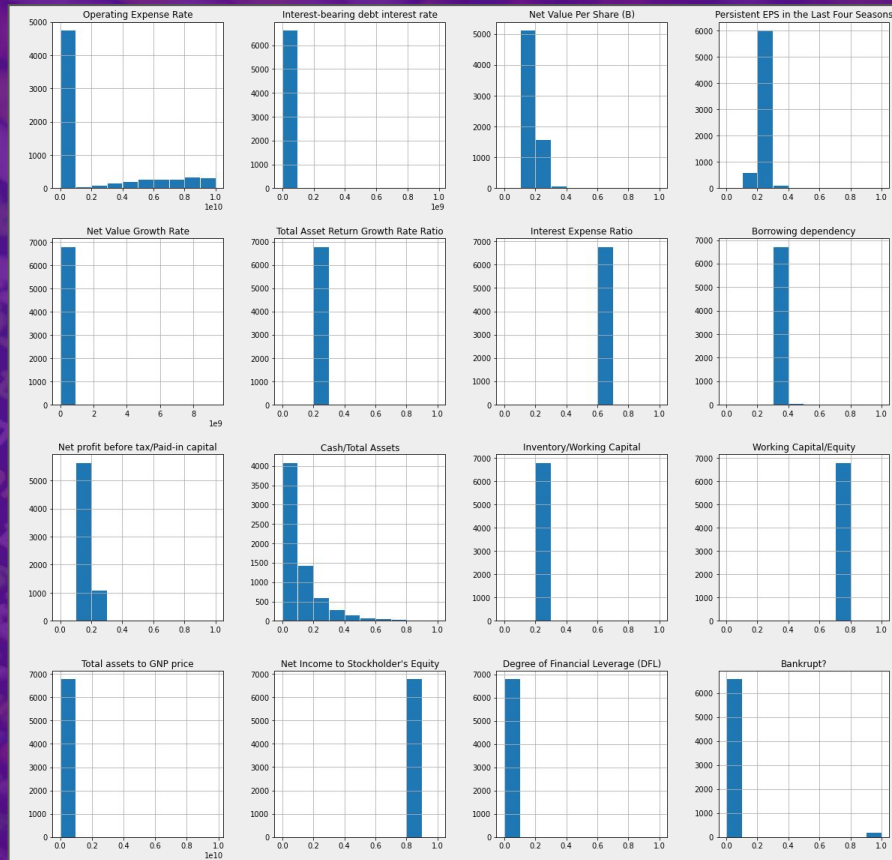| | Operating Expense Rate | Interest-bearing debt interest rate | Net Value Per Share (B) | Persistent EPS in the Last Four Seasons | Net Value Growth Rate | Total Asset Return Growth Rate Ratio | Interest Expense Ratio | Borrowing dependency | Net profit before tax/Paid-in capital | Cash/Total Assets | Inventory/Working Capital | Working Capital/Equity | Total assets to GNP price | Net Income to Stockholder's Equity | Degree of Financial Leverage (DFL) | Bankrupt? |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| count | 6.819000e+03 | 6.819000e+03 | 6819.000000 | 6819.000000 | 6.819000e+03 | 6819.000000 | 6819.000000 | 6819.000000 | 6819.000000 | 6819.000000 | 6819.000000 | 6819.000000 | 6.819000e+03 | 6819.000000 | 6819.000000 | 6819.000000 |
| mean | 1.995347e+09 | 1.644801e+07 | 0.190661 | 0.228813 | 1.566212e+06 | 0.264248 | 0.630991 | 0.374654 | 0.182715 | 0.124095 | 0.277395 | 0.735817 | 1.862942e+07 | 0.840402 | 0.027541 | 0.032263 |
| std | 3.237684e+09 | 1.082750e+08 | 0.033390 | 0.033263 | 1.141594e+08 | 0.009634 | 0.011238 | 0.016286 | 0.030785 | 0.139251 | 0.010469 | 0.011678 | 3.764501e+08 | 0.014523 | 0.015668 | 0.176710 |
| min | 0.000000e+00 | 0.000000e+00 | 0.000000 | 0.000000 | 0.000000e+00 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000e+00 | 0.000000 | 0.000000 | 0.000000 |
| 25% | 1.566874e-04 | 2.030203e-04 | 0.173613 | 0.214711 | 4.409689e-04 | 0.263759 | 0.630612 | 0.370168 | 0.169376 | 0.033543 | 0.277034 | 0.733612 | 9.036205e-04 | 0.840115 | 0.026791 | 0.000000 |
| 50% | 2.777589e-04 | 3.210321e-04 | 0.184400 | 0.224544 | 4.619555e-04 | 0.264050 | 0.630698 | 0.372624 | 0.178456 | 0.074887 | 0.277178 | 0.736013 | 2.085213e-03 | 0.841179 | 0.026808 | 0.000000 |
| 75% | 4.145000e+09 | 5.325533e-04 | 0.199570 | 0.238820 | 4.993621e-04 | 0.264388 | 0.631125 | 0.376271 | 0.191607 | 0.161073 | 0.277429 | 0.738560 | 5.269777e-03 | 0.842357 | 0.026913 | 0.000000 |
| max | 9.990000e+09 | 9.900000e+08 | 1.000000 | 1.000000 | 9.330000e+09 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 9.820000e+09 | 1.000000 | 1.000000 | 1.000000 |

(6819, 16)

# Data

## Looking for Correlations

- Find that top correlations are:
  - *Net profit before tax/Paid-in capital vs. Persistent EPS in the Last Four Seasons*
  - *Persistent EPS in the Last Four Seasons vs. Net Value Per Share (B)*
  - *Net Profit Before Tax/Paid-In Capital vs. Net Value Per Share (B)*

# Data

# Data

## Patterns

- Companies with a low 'Net profit before tax/Paid-in capital', 'Persistent EPS in the Last Four Seasons' and 'Net Value Per Share (B)' tend to go bankrupt

# Data

● Comparing median of each feature for bankrupt vs non-bankrupt companies



| Bankrupt? | Operating Expense Rate | Interest-bearing debt interest rate | Net Value Per Share (B) | Persistent EPS in the Last Four Seasons | Net Value Growth Rate | Total Asset Return Growth Rate Ratio | Interest Expense Ratio | Borrowing dependency | Net profit before tax/Paid-in capital | Cash/Total Assets | Inventory/Working Capital | Working Capital/Equity | Total assets to GNP price | Net Income to Stockholder's Equity | Degree of Financial Leverage (DFL) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.000276 | 0.000317 | 0.185074 | 0.225111 | 0.000463 | 0.264057 | 0.630703 | 0.372474 | 0.179021 | 0.077684 | 0.277179 | 0.736072 | 0.002063 | 0.841232 | 0.026810 |
| 1 | 0.000335 | 0.000499 | 0.158021 | 0.195944 | 0.000396 | 0.263724 | 0.630283 | 0.382655 | 0.154012 | 0.023755 | 0.276981 | 0.732669 | 0.003853 | 0.836707 | 0.026689 |

# Data

## Conclusions of Data Analysis

- Since companies with a low 'Net profit before tax/Paid-in capital', 'Persistent EPS in the Last Four Seasons' and 'Net Value Per Share (B)' are more likely to go bankrupt, we will want to use a KNN model as we have clear clustering
- **Interest-bearing Debt Interest Rate:** When high, tends to bankruptcy
- **Total Assets to GNP Price:** When high, tends to bankruptcy
- **Cash/Total Assets:** When low, tends to bankruptcy

# Models & Evaluation

According to the Data analysis, for this prediction we would train our data through different machine learning algorithms such as:

1. **K-Nearest Neighbors**
2. **Gradient Boosting**

# Models & Evaluation

**K-Nearest Neighbor**

- This model is one of the most commonly used supervised machine learning algorithms.
- It is mostly used in solving classification problems and is based on the nearest neighbor principle.
- This KNN model was trained with the features identified as most predictive during data analysis process: **Net profit before tax/Paid-in capital, Persistent EPS in the Last Four Seasons, Net Value Per Share (B), Interest-bearing debt interest rate, Total assets to GNP price, Cash/Total Assets**
- Evaluated by creating confusion matrix and finding sensitivity scores

# Models & Evaluation

**K-Nearest Neighbor**

```
best number of neighbors: 1
best training set sensitivtiy score : 1.000
best test set sensitivity score: 0.226
```

```
training set score : 0.97
test set score: 0.96
Time:  0.31626288700044825
```

# Models & Evaluation

**Gradient Boosting**

- This gives us prediction model in the method of transformation from weak learners into strong learners.
- First, it creates a starting leaf and then creates new trees by taking into account the errors that occur.
- This process is continued until a better result cannot be obtained.
- Here, we will apply classifier on our reduced data and then on the whole dataset.
- Doing two gradient boosting models and evaluating based on sensitivity scores and time to completion

# Models & Evaluation

## Gradient Boosting

```
training set sensitivity score : 0.30
test set sensitivity score: 0.19
Time:  0.3831532900003367
```

Gradient boosting model performance with reduced (n=6) feature set

```
training set sensitivity score : 0.30
test set sensitivity score: 0.23
Time:  5.036059335999198
```

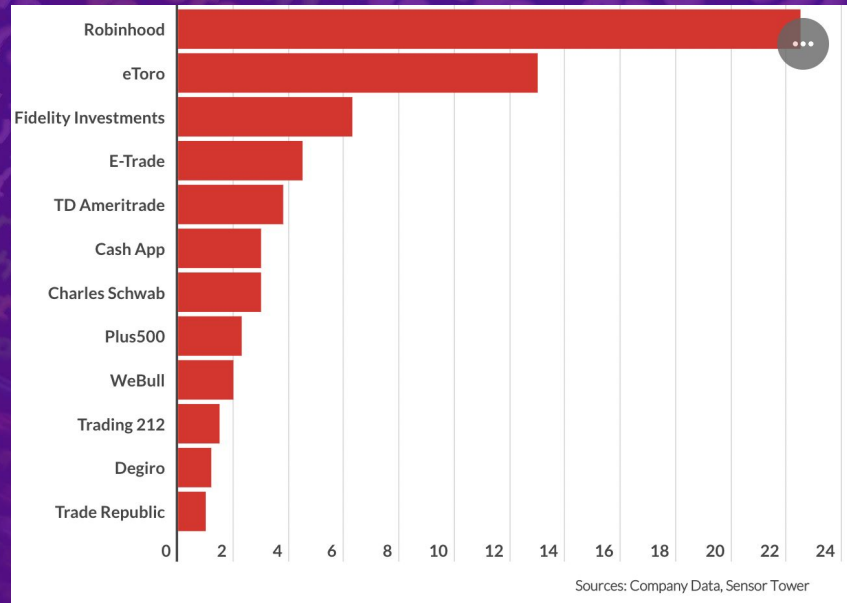Gradient boosting model performance with total original (n=95) feature set

# Deployment

## Business Model

- Our predictive tool will be packaged as a B2B SaaS product
- We'll license the software to consumer-focused stock investment apps (e.g. TD Ameritrade, Robinhood) and charge them based on the number of users

# Deployment

## Stock Trading App Users by App (in millions)



Sources: Company Data, Sensor Tower

NYU | STERN

# Deployment

**Issues and Risks to be Aware of**

- SEC regulations
  - Since our initial target customer will be investment apps, our product has to be in compliance with SEC regulations
- Consumer volatility
  - Investment apps, like Robinhood, tend to encourage emotionally-based reactions instead of logic-based decisions so consumers turn to social media or forums for advice (e.g. Wallstreetbets vs. Robinhood)
- Possible backlash from companies
  - If our tool works as planned, there will be easier access to indicators and information that predict a company's risk of bankruptcy
- Payment-for-order-flow (PFOF) model
  - Many zero-commission investment apps use this process to generate income. The PFOF model bundles trades and sends them to a third-party market maker which then compensates the stockbroker for making the trade. Ultimately, investment apps are selling the users data and monetizing it.

# Deployment

## Conclusion

- As we can see, there are inherent risks with selling financial services products,  specifically when selling directly to consumers
- However, we believe that because we have a B2B model and by aligning ourselves with the SEC's regulations, as well federal and GDRP privacy and security laws, we're able to mitigate some of the associated risks

# Appendix

## Section Authors

- **Business Case:** Brandon Brooks, Nathan Silverglate

- **Data:** Nathan Silverglate

- **Models:** Mayuresh Deolekar, Nathan Silverglate

- **Evaluation:** Mayuresh Deolekar, Nathan Silverglate

- **Deployment:** Brandon Brooks

NYU | STERN

# Appendix

**Sources**

1. https://www.kaggle.com/fedesoriano/company-bankruptcy-prediction
2. Deron Liang, Chia-Chi Lu, Chih-Fong Tsai, Guan-An Shih, Financial ratios and corporate governance indicators in bankruptcy prediction: A comprehensive study, European Journal of Operational Research, Volume 252, Issue 2, 2016, Pages 561-572

NYU | STERN