# Abnormality in Non-Profit Financial Data

Rameen Mahdavi
Github: https://github.com/rmahdavi
Email:Ramahdavi@gmail.com

# Determining the right Non-Profit is hard.

- Historically Limited Data Available

- Little Incentive to release or collect data

- Difficult to measure performance

# But determining the organization not to invest in might be easier.

## Objective

Create an outlier model that will detect higher risk organizations that might need additional scrutiny and see if we can make any inferences.

# Data

- Non-Profit Organization are required to file form 990 with the IRS to keep their tax exempt status
- Recently the IRS released the full form for all 220,000+ nonprofits that filed and posted parsable xmls on Amazon Web Services

# Featurize and Clean

~220,000 Nonprofits

~250 Features

Everything from Real Estate holdings to Indoor tanning expenses

→

~50,000 Nonprofits

Organizations with more 1,000,000 Revenu and 500,000 Donations

With features

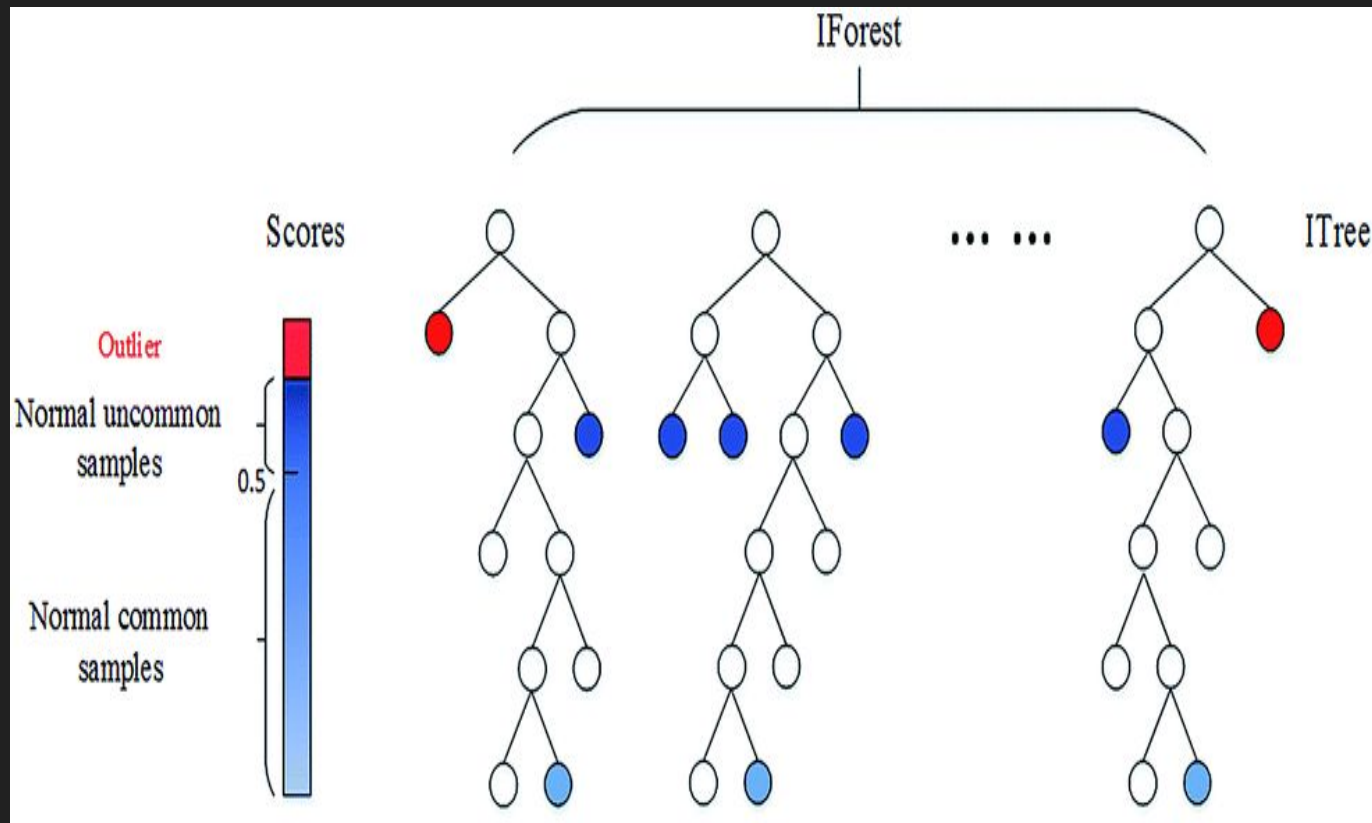| Structure |
| --- |
| 1. Executive Compensation |
| 2.Leverage |
| 3.Solvency |
| **Accounting Manipulation** |
| 4.Deferred Expenses Ratio |
| 5.Deferred Revenues Ratio |
| 6.Depreciation Rate |
| **Performance Metrics** |
| 7.Fundraising Efficiency |
| 8.Surplus Margin |

# Isolation Forest Model

- Good way to identify outliers in Multidimensional data sets

# How it works

- Randomly select a feature and Randomly select a split between the min and max of the feature
- Split until observation is isolated
- Closer the terminal node is root node the more abnormal the observation is
- Perform on a number of trees and average the distances to get an abnormality score

# TSNE Display showing the different clusters of outlier organizations

# What are we detecting?

| Structure |
|---|
| 1. Executive Compensation |
| 2. Leverage |
| 3. Solvency |

| Accounting Manipulation |
|---|
| 4. Deferred Expenses Ratio |
| 5. Deferred Revenues Ratio |
| 6. Depreciation Rate |

| Performance Metrics |
|---|
| 7. Fundraising Efficiency |
| 8. Surplus Margin |

Seems we are identifying non-profits that share similar characteristics and can be classified into distinct clusters

- Low (Assets, Liabilities), High (Revenue and Expenses) cluster - identified by structure metrics
- The Bad Cluster - identified by high judgment metrics

# Thank You

Github: https://github.com/rmahdavi
Email:Ramahdavi@gmail.com

# Appendix:

# OLS Regression on Abnormality Score

```
                            OLS Regression Results
==============================================================================
Dep. Variable:                     AS   R-squared:                       0.241
Model:                            OLS   Adj. R-squared:                  0.241
Method:                 Least Squares   F-statistic:                     2037.
Date:                Tue, 18 Oct 2016   Prob (F-statistic):               0.00
Time:                        23:20:18   Log-Likelihood:                 35193.
No. Observations:               51354   AIC:                         -7.037e+04
Df Residuals:                   51346   BIC:                         -7.030e+04
Df Model:                           8
Covariance Type:            nonrobust
==============================================================================
                 coef    std err          t      P>|t|      [95.0% Conf. Int.]
------------------------------------------------------------------------------
field1       -2.06e-06   4.09e-07     -5.034      0.000    -2.86e-06 -1.26e-06
field2          0.0010      0.000      7.243      0.000        0.001     0.001
field3      -2.896e-07      1e-07     -2.884      0.004    -4.86e-07 -9.28e-08
field4          0.1273      0.004     34.429      0.000        0.120     0.135
field5          0.0446      0.003     14.028      0.000        0.038     0.051
field6          0.2221      0.002     95.015      0.000        0.218     0.227
field7          0.0531      0.001     40.564      0.000        0.051     0.056
field8      -8.589e-08   9.12e-08     -0.942      0.346    -2.65e-07  9.28e-08
==============================================================================
Omnibus:                    28041.573   Durbin-Watson:                   0.926
Prob(Omnibus):                  0.000   Jarque-Bera (JB):         40966739.942
Skew:                           1.074   Prob(JB):                         0.00
Kurtosis:                     141.351   Cond. No.                     4.14e+04
==============================================================================

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 4.14e+04. This might indicate that there are
strong multicollinearity or other numerical problems.
```

## Appendix:
## Random Forest Feature Importance

```
In [101]: log_model.forest_analysis(filepath, 'outlier_IF')
Out[101]:
array([ 0.14680097,  0.14777542,  0.21802028,  0.16047295,  0.0785647 ,
        0.12347648,  0.08946984,  0.03541936])
```