# CMPT 733 Machine learning to detect misstated financial statements
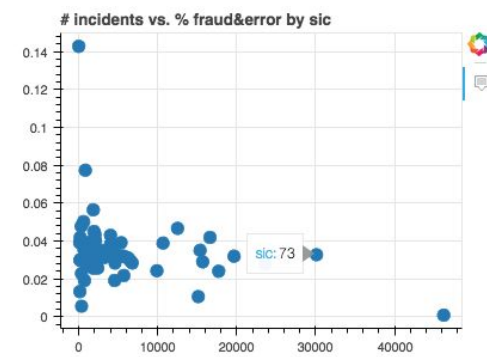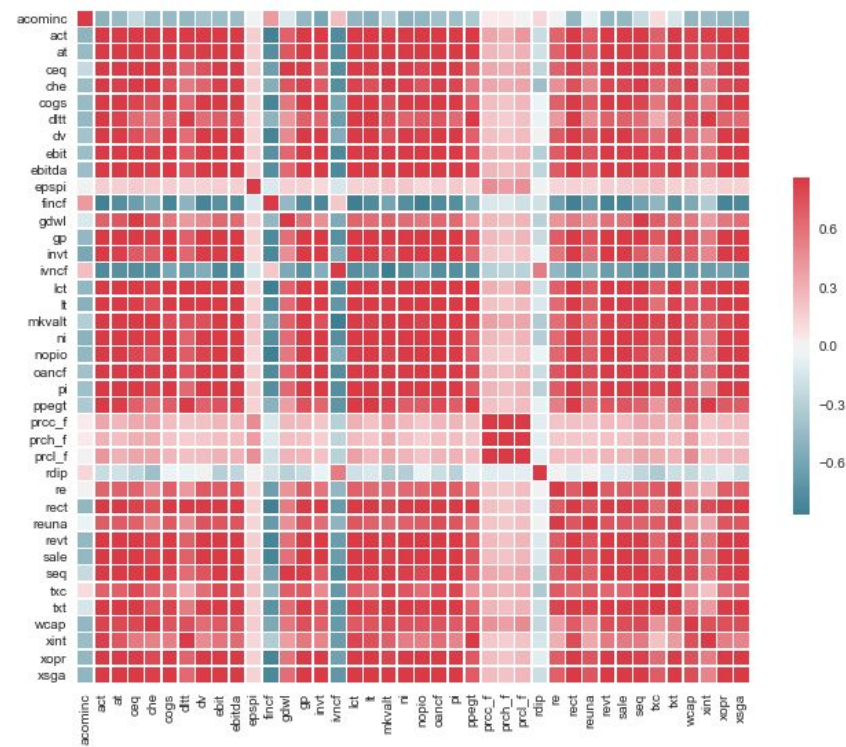
Katrina Ni, Leiling Tao

# Motivation

- Financial statement fraud is deliberate misrepresentation, misstatement or omission of financial statement data, which accounts for 10% of white collar crimes (Association of Certified Fraud Examiners).
- Typically, financial statement contains thousands of fields which follow thousands of rules, presenting a daunting challenge for auditors to detect abnormal fields.
- We aim to automate the process of pre-screening potentially misstated financial statements by using machine learning techniques.
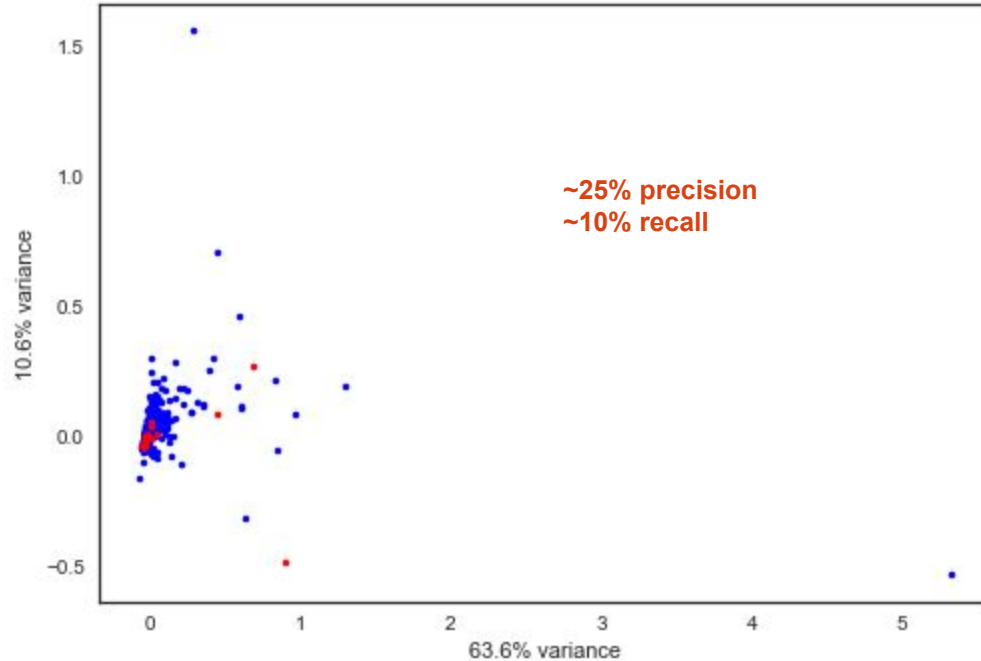
# Progress Report

| Data Collection | Data Preprocessing | EDA |
|---|---|---|
| Data source: WRDS - Wharton Research Data Services<br><br>Raw data ranges from Jan 1980 to Mar 2018 with 435497 observations and 981 features | Exclude observations with missing or zero values for total assets or total equity (~15%)<br><br>Feature selection:<br>• can identify restatement,<br>• have significant effect (control variables such as industry group),<br>• accounting terms (4 levels of importance, starting with the most important ones)<br>drop features with > 50% missing values<br><br>Extract reasons of restatement: <u>fraud and error (~2.7%)</u>, merger and acquisition (~1.4%), change in accounting rules (~0.1%) | Extremely unbalanced classes => unsupervised learning is probably a better option<br><br>Highly correlated accounting variables => Looking for dimensionality reduction<br><br>Look into features that have significant effect according to prior research => Focus on a specific country (USA) or an industry group in a specific year first |

# incidents vs. % fraud&error by fyear

# incidents vs. % fraud&error by sic

Need to focus on a subset of the data with reasonable
number of observations and error/fraud rate
=> start with **USA, year 2004 and sic 73**
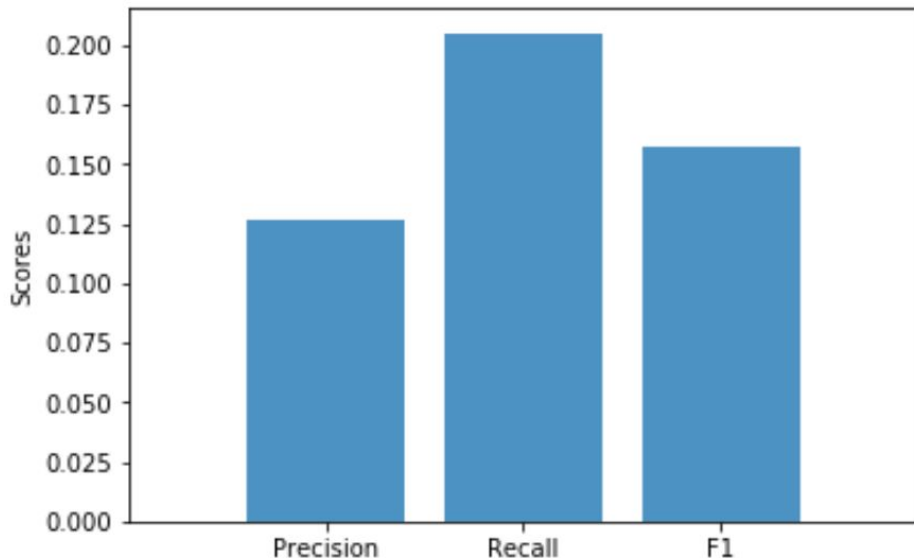
Most accounting features are highly correlated

**Initial Dimensional Reduction with USA, year 2004, industry code 73 (32 abnormalities among 831 data points)**

~25% precision
~10% recall

# Outlier detection using local outlier factor (LOF)

- Using data from 2014, industry code 73, for each variable, we calculated its percent change compared to 2013 (called horizontal comparison in accounting)
- Using these new features, we tried to detect outliers using local outlier factor



Abnormal changes in features are not good indicators of frauds!

Rather, we will focus on abnormal changes in relationships among variables in the future

# Future Work

| Unsupervised Learning | Supervised Learning | Neural network | External Dataset |
|---|---|---|---|
| Deploy several clustering and outlier detection methods:<br><br>• K-nearest-neighbors<br>• Isolation forest<br>• Time series analysis (ARIMA and LSTM) | Classification methods (need to consider imbalance):<br><br>• Random forest<br>• Gradient boosting<br>• SVM | Focus on relationships and patterns among variables, we will construct separate NN models to explore if and how do the structure of fraudulent statements differ from correct ones | Include industry specific indices to explore whether certain changes are more likely to result in frauds, such as<br><br>• Pension data<br>• Credit ratings<br>• Short interest<br>• Executive compensation |