# Using Machine Learning to Find Fraud from Enron Email and Financial Data

## Meixian Chen

## Project overview

The Enron email database consists of over real 600,000 emails from 158 employees, mostly are senior managers of the Enron Corporation. It was acquired by the Federal Energy Regulatory Commission during its investigation after the company's bankrupt, which was the consequence of willful corporate fraud and corruption. The Enron email database is a precious resource for social studies as well as artificial intelligence research. (https://en.wikipedia.org/wiki/Enron_Corpus)

The original data is available in https://www.cs.cmu.edu/~./enron/.

## Explore the dataset

After cleaning and preprocess, the data we have in this project is as following:
It contains 146 points, each point represents a person labeled whether he/she is a person-of-interest (POI).

18 out of 146 persons are POIs.

For each data point, there are financial features, such as  ['salary', 'deferral_payments', 'total_payments', 'loan_advances', 'bonus', 'restricted_stock_deferred', 'deferred_income', 'total_stock_value',

'expenses', 'exercised_stock_options', 'other', 'long_term_incentive', 'restricted_stock', 'director_fees'],

and email features:  ['to_messages', 'email_address', 'from_poi_to_this_person', 'from_messages', 'from_this_person_to_poi', 'shared_receipt_with_poi']

I plot the points by using some combination of features, the distribution of the data looks normal except one point, 'TOTAL', which has a very large "total_payment" values comparing to the other points. It is an outlier and I remove it from the data.

## Features Selection and Optimization

First (manual features selection), I manually eliminate the features which contains too many missing values. They are: deferral_payment, load_adavances, restrict_stock_deferred, and other.

Second (features generation), the percentage of emails from/to/shared a POI shows better how a person related with POIs than the absolute email number. Hence I generate three new features in the following way:

percent_from_poi_to_this_person = from_poi_to_this_person/to_messages
percent_from_this_person_to_poi = from_this_person_yo_poi/from_messages
percent_shared_receipt = shared_receipt_with_poi/to_messages

Now my features_list contains:
['poi','salary','total_payments',  'bonus',
    'deferred_income', 'total_stock_value', 'expenses',
    'exercised_stock_options', 'long_term_incentive',
    'restricted_stock', 'director_fees',
    'percent_from_this_person_to_poi',
    'percent_from_poi_to_this_person', 'percent_shared_receipt']

Third (features scaling), the email features in the features_list range from 0 and 1, while the financial data have a much larger range (up to millions). In order to achieve a good performance of the classification, I apply min_max_scaler to scale all the features to a range of [0,1].

Fourth(intelligent features selection), I use SelectKBest method to choose the best k features. After the fitting the model with the dataset, we have the scores for all the features (the higher score a feature has, the more powerful it tends to be in predicting if a person is a POI.).

salary 18.575703268
total_payments 8.86672153711
bonus 21.0600017075
deferred_income 11.5955476597
total_stock_value 24.4676540475
expenses 6.23420114051
exercised_stock_options 25.0975415287
long_term_incentive 10.0724545294
restricted_stock 9.34670079105
director_fees 2.10765594328
percent_from_this_person_to_poi 16.6417070705
percent_from_poi_to_this_person 3.21076191697
percent_shared_receipt 9.29623187148

Salary, bonus, stock and how often a person write a POI tend to be the most relevant features in predicting a POI.

Fifth (reduce features dimension), I apply PCA to reduce the features dimension in order to achieve better performance of the latter steps.

## Evaluation metric and Validation strategy

I use the scores of Precise and Recall, as well as their harmonic mean F1 to evaluate the classifier method. In this context of this project, Precise is the percentage value of real POIs out of the all POIs the classifier predicts to be; and Recall is the percentage value of real POIs that classifier predicts to be out of all the POIs in the dataset.

For this dataset, I prefer a higher Recall score, which mean a classifier is good to find out suspicious candidates, and then further investigation can be run on the candidates to decide if it is a real POI.

For validation strategy, I use cross validation to evaluate the model. It is important to partition the dataset into training and testing parts. We fit the model with training data and in order to fairly evaluate a model, we use the test data, which the model hasn't seen, to check its predicting ability. Without separating training and testing data, a model tends to tell a better performance than what it can actually do on new data.

I use StratifiedShuffleSplit cross-validation method, which is to randomly separate the data into training and testing part several times, and calculate the average evaluation scores.

## The First try of different classifiers

I try the following classifier: GaussianNB, AdaBoostClassifier, SVC, NearestCentroid, and RandomForestClassifier.

Without tuning the parameter, their performances are as following:

GaussianNB:
Accuracy: 0.86000   Precision: 0.47059   Recall: 0.40000      F1: 0.43243

AdaBoostClassifier

Accuracy: 0.84000   Precision: 0.35714   Recall: 0.25000      F1: 0.29412

SVC
Precision or recall may be undefined due to a lack of true positive predicitons
(And SVC is very hard to tuning to get a reasonable performance for this dataset)

NearestCentroid
Accuracy: 0.86667   Precision: 0.50000   Recall: 0.35000      F1: 0.41176

RandomForestClassifier:
Accuracy: 0.86667   Precision: 0.50000   Recall: 0.10000      F1: 0.16667


Among the above classifier, AdaBoostClassifier and
RandomForestClassifier are very slow in passing the test from tester.py.

I decide to tuning the parameters of the algorithm which use
GaussianNB and NearestCentroid.


# Tune the algorithm

I use GridSearchCV to automatically explore different combination of
values from the parameter grid, and use the one that gain the highest
"f1" score in the cross validation.

The parameter grid contains different k values for SelectKBest method
and n_components values for PCA.

For GaussianNB, the best parameters are SelectKBest_k = 10 and
pca_n_components = 5.
The performance evaluated by tester.py is
Accuracy: 0.85333   Precision: 0.43750   Recall: 0.35000      F1: 0.38889

For NearestCentroid, the best parameters are SelectKBest_k = 6 and
pca_n_components = 3.
The performance evaluated by tester.py is

Accuracy: 0.84667     Precision: 0.43478     Recall: 0.50000   F1: 0.46512
(Both Precision and Recall are higher than the requirement (0.3) for this project)

Therefore, my classifier looks like this:
clf = Pipeline( [('scaler', min_max_scaler),('SKB', SelectKBest(k=6)),('reduce_dim', PCA(n_components = 3)), ('clf', NearestCentroid )])