# Project Overview

In 2000, Enron was one of the largest companies in the United States. By 2002, it had collapsed into bankruptcy due to widespread corporate fraud. In the resulting Federal investigation, a significant amount of typically confidential information entered into the public record, including tens of thousands of emails and detailed financial data for top executives. In this project, I will play detective, and put my new skills to use by building a person of interest identifier based on financial and email data made public as a result of the Enron scandal. To assist me in my detective work, the data combined with a hand-generated list of persons of interest in the fraud case, which means individuals who were indicted, reached a settlement or plea deal with the government, or testified in exchange for prosecution immunity.

# Summary

The goal of this project is to develop an algorithm that could have identified people of interest related to its eventual financial disaster.  This project will use machine learning with supervised learning since we know the outcome of guilty parties.  The data used in this project is comprised of many financials and emails from dozens of former Enron employees.  Financial data includes categories such as salary, bonus, exercised stock options, etc.  The emails data is all the available email traffic to or from a former employee.

Two outliers were identified and removed: *TOTAL* and *THE TRAVEL AGENCY IN THE PARK (TAP)*.  Both were removed mainly because they were not people.  TOTAL is obviously not relevant for determining specific people engaged in financial misconduct.  TAP was removed as we are not focused on potential companies of financial misconduct—only people.  TAP was 50% owned by Sharon Lay, Ken Lay's sister.  TAP received commissions for managing travel books for Enron.[1]  While the family connection may be suspicious, it is outside the scope of the project.

# Features

All available features from the data set were used except the actual email addresses.  Two new features were created but not used: total_comp and fraction_to_poi.  The feature total_comp sums salary, bonus, total_stock_value, and exercised_stock_options.  Fraction_to_poi generates the fraction of emails to and from the POI.  The idea was to compress features into a bigger metric; however, neither feature seemed to help with accuracy, precision, or recall.

The features list was pass through SelectKBest to identify the top 6 features.  There were a total of 18 identified features.  Every number of top features from 4 to 12 was tested.  A feature quantity of 6 produced the highest accuracy, precision, and recall.  Here are the top 6 features and scores using SelectKBest:

```
#1 exercised_stock_options with score 24.8150797332
#2 total_stock_value with score 24.1828986786
#3 bonus with score 20.7922520472
#4 salary with score 18.2896840434
#5 deferred_income with score 11.4584765793
#6 long_term_incentive with score 9.92218601319
```

The feature loan_advances have the highest number of *NaN* values.  Below is a list of Total NaN's by feature in the data set:

---

[1] Texas agency takes a huge hit from Enron's fall. March 5, 2002.  http://www.travelweekly.com/Travel-News/Travel-Agent-Issues/Texas-agency-takes-a-huge-hit-from-Enron-s-fall

```
* salary with 50 NaN's
* to_messages with 58 NaN's
* deferral_payments with 106 NaN's
* total_payments with 21 NaN's
* exercised_stock_options with 43 NaN's
* bonus with 63 NaN's
* director_fees with 128 NaN's
* restricted_stock_deferred with 127 NaN's
* total_stock_value with 19 NaN's
* expenses with 50 NaN's
* from_poi_to_this_person with 58 NaN's
* loan_advances with 141 NaN's
* from_messages with 58 NaN's
* from_this_person_to_poi with 58 NaN's
* deferred_income with 96 NaN's
* shared_receipt_with_poi with 58 NaN's
* restricted_stock with 35 NaN's
* long_term_incentive with 79 NaN's
```

## Algorithm

Two algorithms were tested with the dataset: Guassian Naïve Bayes and SVC.  Scaling was not utilized with GNB as it was not required.  Naïve Bayes resulted the following:

- Accuracy score:  0.928571428571
- Precision:  0.5
- Recall:  0.666666666667

SVC resulted in the following when used with MinMaxScaler for scaling:

- Accuracy score:  0.928571428571
- Precision:  0.5
- Recall:  0.333333333333

The results where almost identical; however, GNB had a higher recall.

## Parameter Tuning

Parameter tuning can be a very important process when algorithms require it.  In this project, Guassian Naïve Bayes was used which does not require parameter tuning.  SVC does utilize parameter tuning which was executed in this project.  The kernel and C parameters were investigated for SVC.  In simple terms, the lower the C value is the higher the probability of misclassification.  Therefore, we experimented with different C values before settling on 1,000.

## Validation

Validation is an important process with machine learning.  Without proper validation we can create scenarios of overfitting amongst other potential issues.  The focus of this project was to get the accuracy, precision, and recall as close to a value of 1 as possible.  The higher the precision and recall, the lower the likelihood our models are subject to overfitting.

## Metrics

Precision and recall are used to gauge the efficacy of the investigated algorithms.  Precision is the ratio of true positives to the sum of true positives and false positives.  Recall is the ratio of true positives to

the sum of true positives and false negatives.[2]  The closer to 1 for both metrics, the lower the error rate.  An example related to this project would be if an individual is falsely identified as a person of interested when in fact they are not (i.e. precision).  Another example would be if an individual is not identified as a person of interest when in fact they are (i.e. recall).  These examples are the reason we want to make sure all persons of interest are identified and no one is falsely identified.

---

[2] Precision and Recall. September 16, 2016. https://en.wikipedia.org/wiki/Precision_and_recall