

1. The purpose of this project is to determine “persons of interest” (POI) in the infamous corporate scandal that led to Enron’s demise in the early 2000’s. I used publically available data of 145 Enron employees to train a machine learning algorithm to detect POI. Each employee has 21 features pertaining to their email metadata and financial information (salary, bonus, stock value, etc.). Of the 144 employees in the dataset, 18 are listed as POI. There is a glaring outlier in this dataset with the name “TOTAL” that includes the sum of the financial statistics of all 144 employees. Since this observation is not relevant to my analysis, I removed it from the dataset. I also removed all values that had 18 or more NaN features. There were 5 values that met this criteria:

- GRAMM WENDY L
- LOCKHART EUGENE E
- THE TRAVEL AGENCY IN THE PARK
- WHALEY DAVID A
- WROBEL BRUCE

After removing these values, the total amount of NaN values dropped from 1358 to 1260. I felt that these values were not providing much information to the classifier and could’ve been hurting its performance. Machine learning is useful in this application because it learns the patterns in the data and is able to associate them with a label (POI or non-POI). These processes are known as “training” and “classifying”, respectively. Once a classifier is trained, it has the capability to predict labels for which there is no label already assigned.

2. In the first step of the feature selection process, I determined the Pearson correlation coefficient for all features in the dataset, pairwise. For any pairs that had a coefficient value larger than 0.75, I removed one from the pair. This was done to reduce the possibility of over-fitting. Next, I used the “Select K Best” algorithm to reduce features that didn’t contribute to maximizing the f1 score. I did not scale the features because feature distance is not important in the algorithm that I used (AdaBoost with decision tree learning). I engineered and tested three additional features – percentage of incoming emails from POI’s, percentage of outgoing emails to POI’s, and percentage of incoming emails that share receipt with a POI. I created these features because I felt that an email sum was susceptible to bias toward those who generally sent more emails or had a longer tenure with the company. I thought that an employee’s percentage of emails with POI’s, rather than the count, would reveal more about their affiliation with POI’s. In the end, I did not include these additional features because after testing, they did not improve the performance of my classifier. These are the features that I ended up using, along with their scores in the “Select K Best” algorithm:

- salary: 17.4320873825
- deferral_payments: 0.24851659019
- total_payments: 8.47223400901
- bonus: 19.9898755267
- restricted_stock_deferred: 0.0670959907812

- deferred_income: 11.0477590807
- total_stock_value: 23.3373703767
- expenses: 5.67653901906
- long_term_incentive: 9.47251713657
- director_fees: 2.00501730869
- from_poi_to_this_person: 4.94057196272
- from_messages: 0.187200230841
- from_this_person_to_poi: 2.25165650525
- shared_receipt_with_poi: 8.1199380825

3. I used the AdaBoost algorithm with Decision Tree Learning as my weak learner. Using Grid Search CV as a validator, it got a “best score” of 0.345. I also tried the Naïve Bayes Classifier, which gave me a “best score” of 0.309. Since AdaBoost had a better score, I decided to use it for the final algorithm.
4. Parameter tuning is the process by which parameter inputs to a classifier are altered to improve the performance of an algorithm. If this process isn’t done properly, it could lead to bias, over-fitting, and ultimately a poor classifier. I tuned my parameters using “Grid Search Cross Validation”. I tuned both the “n_estimators” and “learning_rate” parameters of the AdaBoost Classifier, and scored using an f1 score. Grid Search found that the best parameters were `n_estimators = 75` and `learning_rate = 0.9`. I also tuned the “k” value of Select K Best, and found that `k = ‘all’` to be best.
5. Validation is the process of measuring the performance of a classifier. It consists of splitting data into training and testing sets, and determining how well a classifier is able to predict the labels of the test set, after being fit on a training set. A mistake in validation occurs when the same data is used to train and test a classifier. This leads to bias when evaluating its performance because the classifier is already really good at classifying the data it was trained on. I evaluated my analysis twice, first using a “Grid Search CV” to decide between two algorithms, and later using a “Stratified Shuffle Split” to measure more specific scores such as precision and recall. This algorithm splits the dataset into several folds, while preserving the percentage of samples in each class. Therefore, each fold will have the same proportion of POI’s to non-POI’s. Note: The “Stratified Shuffle Split” was provided in `tester.py`.
6. Using “Stratified Shuffle Split”, I got a precision value of 0.47964 and a recall value of 0.31800. Precision is defined as the ratio of correctly identified true values and the total predicted true values. The predicted true values contain both incorrect true values (false positives) and correct true values (true positives). Out of 1400 predictions, the classifier predicted 636 true values correctly. Recall is defined as the ratio of correctly identified true values and the actual true values. In this case, on average roughly 8 of 18 POI’s were predicted correctly.