

Identifying Fraud at Enron using Emails and Financial Data

Udacity Data Analysis Nanodegree – Project 4

Introduction

Enron Corporation was one of the world's largest energy companies before it went bankrupt in December 2001. It was revealed that the bankruptcy was caused by 'institutionalized, systematic and creatively planned accounting fraud'. Subsequent legal trials found several members of Enron's board guilty of being involved in the financial fraud.

For this project I have attempted to identify the Persons of Interest (POI) involved in fraudulent activity which contributed to the collapse of Enron. To do this I have used several machine learning techniques to help characterise important features within the financial and email dataset and then used these to help predict which of the employees might be POIs.

Analysis

The dataset was made up of financial and email features. The email data contained both the content of the email messages and also associated metadata (sender and recipient information), the financial data had a range of features including salary, stock options, bonuses, expenses etc. There are 146 entries in the dataset with each entry having 21 features. It is worth noting that not every feature had a completed value and that the dataset contained a subset of the total Enron employees.

The original dataset contained an entry for the **TOTAL** of all the financial features of each of the Enron employees. This was obvious when plotting the data as it was several magnitudes larger than any other data point (Fig 1), I therefore removed this line from the dataset. Once this line was removed then the same plot (Fig 2) showed only a handful of employees as outliers, however this was expected in a large organisation.

Fig 1.

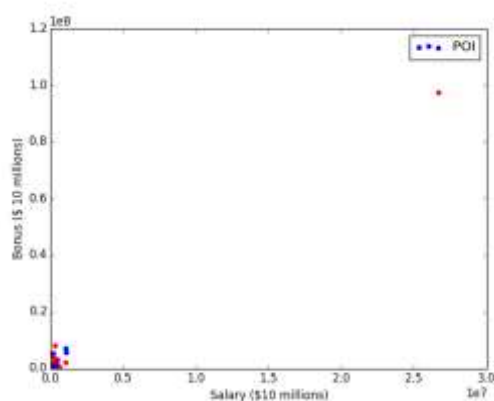
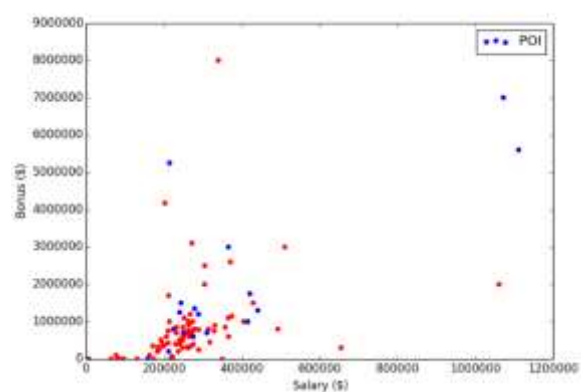


Fig 2.



I also removed the entry for **Eugene E. Lockhart** since this row had no values for any of the features. Finally after looking at the data in more detail I removed an entry for **The Travel Agency in the Park** as I was only interested in people for my machine learning model. After removing these three entries I was left with 143 records in the dataset.

Feature Selection

In order to optimise and select the most relevant features I used the scikit-learn module **SelectKBest** to identify the top ten most influential features.

Feature	Score
exercised_stock_options	24.815
total_stock_value	24.183
Bonus	20.792
Salary	18.290
fraction_emails_to_POI	16.410
deferred_income	11.458
long_term_incentive	9.922
restricted_stock	9.213
total_payments	8.773
shared_receipt_with_POI	8.589

When I first ran the **SelectKBest** algorithm the only email feature in the top ten was `shared_receipt_with_POI`. I was surprised that neither `from_this_person_to_POI` nor `from_POI_to_this_person` were included in the most influential features as I would have thought the level of interaction with a POI would have been a strong indicator of also being a POI. Plotting the data (Fig 3) it didn't really show a strong correlation, hence the low score. However when I plotted the fraction of the total emails that went to a POI or came from a POI (Fig 4) I found a much clearer link. This made sense as email volumes can vary quite dramatically and so a percentage of the total is a far better indicator. Once I re-ran the **SelectKBest** the `fraction_emails_to_POI` became the 5th strongest feature. I used the above ten most important features in my final analysis.

Fig 3.

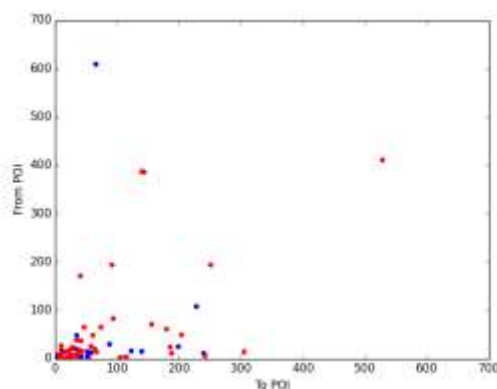
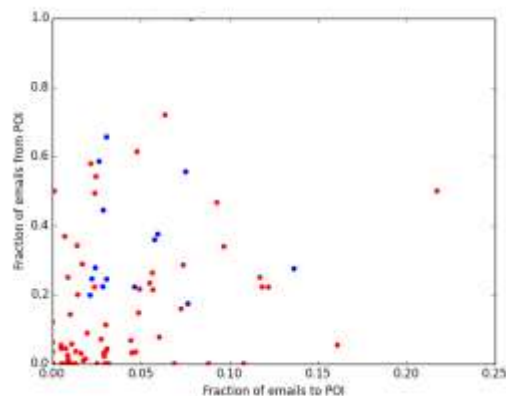


Fig 4.



Algorithm

I looked at four different machine learning classifiers; Naïve Bayes, Decision Tree, K Means Clustering and Support Vector Machines (SVM).

Prior to training the K Means Clustering, Decision Tree and SVM classifiers I scaled all the features using the Min-Max feature scaler. This was incredibly important as the features had different units of measurement which varied by several orders of magnitude. By scaling the features it meant that I could use these classifiers and know that the features would be weighted evenly.

Algorithm	Accuracy	Precision	Recall
Naïve Bayes	0.33547	0.14750	0.83350
Decision Tree	0.82047	0.32632	0.32550
K Means Clustering	0.83760	0.23086	0.09350
SVM	0.84119	0.26882	0.15700

The tuning parameters and the detailed scoring for each model can be found below:

1) **Naïve Bayes** (no scaling or tuning):

Accuracy: 0.33547, Precision: 0.14750, Recall: 0.83350, F1: 0.25064, F2: 0.43182

Total predictions: 15000, True positives: 1667, False positives: 9635, False negatives: 333, True negatives: 3365

2) **Decision Tree** (features scaled using Min-Max scaler):

GridSearchCV(cv=None, estimator=DecisionTreeClassifier(compute_importances=None, criterion=gini, max_depth=None, max_features=None, min_density=None, min_samples_leaf=1, min_samples_split=2, random_state=None, splitter=best)

Accuracy: 0.82047, Precision: 0.32632, Recall: 0.32550, F1: 0.32591, F2: 0.32566

Total predictions: 15000, True positives: 651, False positives: 1344, False negatives: 1349, True negatives: 11656

3) **K means Clustering** (features scaled using Min-Max scaler)

KMeans(copy_x=True, init=k-means++, max_iter=300, n_clusters=2, n_init=10, n_jobs=1, precompute_distances=True, random_state=None, tol=0.001, verbose=0)

Accuracy: 0.83760, Precision: 0.23086, Recall: 0.09350, F1: 0.13310, F2: 0.10613

Total predictions: 15000, True positives: 187, False positives: 623, False negatives: 1813, True negatives: 12377

4) **Support Vector Machines** (features scaled using Min-Max scaler)

SVC(C=1000, cache_size=200, class_weight=None, coef0=0.0, degree=3, gamma=0.0, kernel=rbf, max_iter=-1, probability=False, random_state=None, shrinking=True, tol=0.001, verbose=False)

Accuracy = 0.84119, Precision = 0.26882, Recall = 0.15700

After testing these different algorithms I decided to use the Decision Tree classifier as it performed best across the three evaluation metrics. I tuned the DecisionTreeClassifier using GridSearchCV. I tuned on both **criterion** (the function used to measure the quality of the split) as well as the **splitter** (the strategy used to choose the split at each node). The best parameters were:

`DecisionTreeClassifier() - {'splitter': 'random', 'criterion': 'gini'}`

Validation

Validation is performed to ensure that a machine learning algorithm generalises well. A classic mistake is overfitting, this is where a model is trained and performs very well on the training dataset but a lot worse on the test datasets. To avoid over-fitting I used **cross_validation** to split the Enron data into a test and a training set. I took the average precision and recall over 1000 randomised tests splitting the data into 30% test and 70% training.

Evaluation

Precision is the rate at which the algorithm correctly predicts the POI. It is the ratio of true positives to the records that are actually POIs. It is calculated by: $\text{True Positives} / (\text{True Positives} + \text{False Negatives})$. My DecisionTreeClassifier had a precision of **0.32632**

Recall measures the sensitivity and refers to the proportion of the POI the model can detect of all the POI. Recall is calculated as $\text{True Positives} / (\text{True Positives} + \text{False Positives})$. My model achieved a recall score of **0.32550**.

In this investigation accuracy was not a good metric given the high number of non-POIs in the dataset. For example if non-POI had been predicted for all records then an accuracy of over 87% could have been achieved.

Conclusion

The Enron dataset provided a very interesting challenge of how to apply machine learning methods to a small but complex dataset. The data needed to be cleaned and then a combination of machine learning classifiers, parameter tuning and cross-validation techniques led to the creation of a reliable predictive model for POIs which tried to maximise the recall and precision while maintaining a high accuracy score.

I enjoyed the challenge of studying machine learning techniques and the data workflows required for data analysis in Python.

References

Udacity module page - <https://www.udacity.com/course/intro-to-machine-learning--ud120>

Enron Wikipedia - <https://en.wikipedia.org/wiki/Enron>

Sci-kit learn page - <http://scikit-learn.org/stable/>

SelectKBest page - http://scikit-learn.org/stable/modules/generated/sklearn.feature_selection.SelectKBest.html

Mix Max Scaler - <http://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.MinMaxScaler.html>

Naïve Bayes - http://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.GaussianNB.html

Decision Trees - <http://scikit-learn.org/stable/modules/tree.html>

K Means Clustering - <http://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html>

Support Vector Machines (SVM) - <http://scikit-learn.org/stable/modules/svm.html>