

Identifying Fraud from Enron Emails and Financial Data

Introduction

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, there was a significant amount of typically confidential information entered into the public record, including tens of thousands of emails and detailed financial data for to executives.

Utilizing the classifiers and techniques taught in Into to Machine Learning Class, I built a classifier to detect if a person is culpable or not.

Short Questions

Summarize for us the goal of this project and how machine learning is useful in trying to accomplish it. As part of your answer, give some background on the dataset and how it can be used to answer the project question. Were there any outliers in the data when you got it, and how did you handle those?

The features in the data fall into three major types, namely financial features, email features and POI labels.

- financial features: 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
- email features: to_messages, email_address, from_poi_to_this_person, from_messages, from_this_person_to_poi, shared_receipt_with_poi
- poi label(A total of 18 entries was labelled as POI)

The goal of this project was simply to leverage the features above in order to mark an individual as a person of interest.

In order to detect outlier and explore the data, I first built a code to convert the dictionary to csv file and is available with the submission. Later, I explored all the values to find out which features had a lot of values and which did not. Immediately looking at the csv, I noticed that the number of data is only 146 and so individually scanning each name - I found the following data points redundant

- Total: Does not convey information pertaining to any individual, hence marked as outlier.
- THE TRAVEL AGENCY IN THE PARK: It does not represent any individual and hence was removed.
- LOCKHART EUGENE E: This record contained no useful data.

After cleaning the data only 143 records remained.

What features did you end up using in your POI identifier, and what selection process did you use to pick them? Did you have to do any scaling? Why or why not? As part of the assignment, you should attempt to engineer your own feature that doesn't come ready-made in the dataset—explain what feature you tried to make, and the rationale behind it. If you used an algorithm like a decision tree, please also give the feature importances of the features that you use.

To pick best features, I used the select k best features of scikit learn. After trying different K values, The K-best approach is an automated univariate feature selection algorithm, and in using it. I decided to go with two algorithms namely K-means clustering and Logistic Regression. In, both cases I had to use different number of features to get the required result

```

'salary': 18.289684043404513
'email_interaction': 5.3993702880944232
'total_payments': 5.3993702880944232
'loan_advances': 8.7727777300916792
'bonus': 7.1840556582887247
'financial_aggregate': 20.792252047181535
'total_stock_value': 15.971747347749965
'shared_receipt_with_poi': 24.182898678566879
'from_poi_to_this_person': 8.589420731682381
'exercised_stock_options': 8.589420731682381
'deferred_income': 5.2434497133749582
'expenses': 24.815079733218194
'restricted_stock': 11.458476579280369
'long_term_incentive': 6.0941733106389453
'fraction_from_poi_email': 9.2128106219771002
'fraction_to_poi_email': 9.9221860131898225

```

For K-means clustering I used a total of 8 features

```

'salary': 18.289684043404513
'bonus': 20.792252047181535
'financial_aggregate': 15.971747347749965
'total_stock_value': 24.182898678566879
'exercised_stock_options': 24.815079733218194
'deferred_income': 11.458476579280369
'fraction_from_poi_email': 9.2128106219771002
'fraction_to_poi_email': 9.9221860131898225

```

The feature:email_interaction did decreased both precision to recall considerably in the case of K-means clustering and hence it wasn't used.

From the above features, I added a total of 4 new features which increased both the recall and precision of Logistic Regression. They are as follows:

- email_interaction_ratio: which was a ratio of the total number of emails to and from a POI to the total emails sent or received.
- financial_aggregate: This was the combined sum of exercised_stock_options, salary, and total_stock_value. This captured both liquid and semi solid wealth an individual has.
- fraction_from_poi_mail: Ratio of mails received from poi to total received mail
- fraction_to_poi_mail: Ratio of mails sent to poi to the total mails sent

justifition

After testing the k value for financial aggregate I got the following: 15.971747347749965 Thus it is one of the significant features contributing to precision and recall. similarly with the new features i.e

- 'fraction_from_poi_email': 9.2128106219771002
- 'fraction_to_poi_email': 9.9221860131898225

Recall for Logistic Regression increased above 0.3.

Usual Features(feature selection was done till the best precision and recall values were obtained):

Classifier	Precision	Recall	Features
Logistic Regression	0.367	0.224	12
K-means Clustering, K=2	0.427	0.307	5

New Features:

Classifier	Precision	Recall	Features
Logistic Regression	0.317	0.300	16
K-means Clustering, K=2	0.340	0.374	8

I scaled all features using a min-max scaler. This ensures features are evenly balanced and overcomes the disparity due to the units of financial and email features.

What algorithm did you end up using? What other one(s) did you try?

After having performed various ml related projects during my undergraduate studies, a two class problem usually is best for logistic regression and K-means clustering. K-means clustering with PCA and mahalanobis distance provides a very fortified technique for two class detection. However, I went with logistic regression as my final algorithm.

I tried several algorithms, with a K-means clustering algorithm performing reasonably sufficient. I also tested a support vector machine, a random forest classifier, and stochastic gradient descent. The best results I got were from logistic regressor.

following were the parameters I tuned: - Logistic regression: C (inverse regularization parameter), tol (tolerance), and class_weight (over/undersampling) - K-means clustering: tol

K-means clustering was initialized with K (n_clusters) of 2 to represent POI and non-POI clusters. It performed well with the 11 set of features.

Auto-weighting in the case of logistic regression caused a dip in precision and hence I decided to use evenly balanced features.

The other algorithms were tuned experimentally, with unremarkable improvement.

What is validation, and what's a classic mistake you can make if you do it wrong? How did you validate your analysis?

Validation is performed to ensure that a machine learning algorithm generalizes well. The classic problem that can occur is over-fitting. This happens when we overfit the training data and perform really well in it due to which there is a considerable dip in the performance in the other two datasets (cross validation and testing dataset).

I validated my result using two techniques - bootstrapping (cleaner.py) - k fold method (run mytester to get the results from k fold method)

Give at least 2 evaluation metrics, and your average performance for each of them. Explain an interpretation of your metrics that says something human-understandable about your algorithm's performance.

I am using precision and recall as the evaluation metric. - Precision: The precision can be interpreted as the likelihood that a person who is identified as a POI is actually a true POI - Recall: Recall measures how likely it is that, given that there's a POI in the test set, this identifier would flag him or her

I did not choose accuracy because with such a small dataset using it would mean that a simple heuristic with all data marked false would give an accuracy of more than 85%

Validation 1 (Stratified K-folds, K=3) Run mytester for this

Classifier	Precision	Recall	Features
Logistic Regression	0.394	0.320	16
K-means Clustering, K=2	0.372	0.363	8

Validation 2 randomised sampling(n=10000)

Classifier	Precision	Recall	Features
Logistic Regression	0.317	0.300	16
K-means Clustering, K=2	0.340	0.374	8

Both algorithms do well inspite of the dataset being noisy. Let us understand what the values are actually speaking to us. Now, in terms of precision and recall you would want a high recall in such a case. Simply, because you want to be suspecting people. Having a high precision would mean that we are looking for too strict of a conditions to flag an individual. Hence, recall plays a high role in such a case.

Conclusion

The dataset was sparse and most algorithms will perform well only if given a decent number of data to learn. Had there been a large dataset random forest would also work very well. However, I am more interested in what a classic anomaly detection system would do in such a case. A simple combination of - K-means - PCA - Mahalanobis distance measure would result in a very robust outlier detection algorithm and probably would have both a high recall and precision value.

Resources and References

- [Introduction to Machine Learning \(Udacity\)](#)
- [Machine Learning \(Stanford/Coursera\)](#)
- [scikit-learn Documentation](#)