

Report

October 17, 2016

0.1 Identifying Fraud from Enron Emails and Financial Data

0.1.1 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 guilty or not. The persons who are guilty are termed POI(Person of Interest) in the context of this problem.

0.1.2 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 goal of this project is to identify fraud from email and financial data available about Enron employees. We are trying to identify each individual as person of interest or not. The dataset consisted of the following features: - **email features:** to_messages, email_address, from_poi_to_this_person, from_messages, from_this_person_to_poi, shared_receipt_with_poi - **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

Apart from the above features, I constructed some binary features that captured the information for missing values in the features. Also, I created a new feature that measured the interaction of an individual with the POIs. The decision to create binary features was taken while exploring the data and finding out that a lot of values for the features were missing. The data exploration was performed on a csv file generated from the data-dict supplied for the project. During exploration, the following two outliers were identified:

- **TOTAL** - This was the total of all the people in the dataset.
- **THE TRAVEL AGENCY IN THE PARK** - This does not belong to any of the employees.

There were only 146 records in total and after removal of above two outliers, only 144 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 does not come ready-made in the dataset – explain what feature you tried to make, and the rationale behind it.

I am using an L1 penalty based logistic regression as a POI Identifier which is also a model for feature selection. It automatically produces sparse models. The following features were not included in this sparse model:

- missing_bonus
- missing_total_stock_value
- missing_total_payments
- missing_to_messages
- missing_loan_advances
- missing_from_this_person_to_poi
- missing_salary
- missing_from_poi_to_this_person
- missing_restricted_stock
- interaction_with_poi

As discussed in the previous question, I created binary features for all features available in the data to capture the missing value information for these features. If the value of a feature was missing, the binary feature was populated with 1 and otherwise with 0. Also, a feature measuring the interaction of an individual with poi was created. However, none of these engineered features proved to be very significant. The top 10 features selected from the random forest classifier are listed below with their respective importance:

Feature	Imp
exercised_stock_options	0.16
bonus	0.14
total_stock_value	0.13
total_payments	0.126
expenses	0.117
salary	0.067
from_poi_to_this_person	0.063
missing_deferred_income	0.053
from_this_person_to_poi	0.050
from_messages	0.030

No feature scaling was performed as tree based models and logistic regression are not impacted by different scales of data.

What algorithm did you end up using? What other one(s) did you try? How did model performance differ between algorithms?

I tried 3 different algorithms for my POI Identifier. I used decision trees, Random Forest and Logistic Regression. However, I chose Logistic regression with L1 penalty for the final model.

The choice for Logistic regression came after validating the model with 10-fold cross validation and tuning the hyperparameters which resulted in a good enough precision and recall for the positive class (POI in this case). The decision tree had a high variance among the results. So, it was dropped and random forest was tested. However, random forest also did not give the desired precision and recall and hence, was discarded.

What does it mean to tune the parameters of an algorithm, and what can happen if you don't do this well? How did you tune the parameters of your particular algorithm?

It is very important to tune the parameters of an algorithm because tuning the parameters helps the algorithm to perform well and generalize well. In order to tune the hyperparameters for all the three models, I used the following approach: - I performed a 80-20 train-test split and kept the testing set aside. - I performed 5 fold cross validation for Decision Tree Classifier while searching for the best parameters using Grid Search. After the best model was identified based on 5-fold cross validation, I tested the model on the final test set. For the other two models, i.e. random forest and logistic regression, I performed 10-fold cross-validation and performed a final test on the test dataset with the best identified model.

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 the machine learning model generalizes well. The most common issue that arises is overfitting wherein a model simply remembers the training set as much as possible and then it is not able to generalize to unseen cases well. When this happens, the performance of the algorithm drops drastically. I used K-fold cross validation and learning and validation curves to validate my analysis (all analysis present in the notebook, `model_selection.ipynb`).

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.

In order to tune the algorithms, I used F1 score as the scoring function as it captures the precision and recall information in a single number. I wanted both my precision and recall to be higher; however, recall is much more important as it will help us tag POIs. We want to flag people who might be fraudulent and then pursue an investigation. In the context of our project, precision measures the accuracy of our POI identifier in identifying the POIs. It gives us the probability of an individual being POI if he has been identified as such by our model. On the other hand, recall measures the ability of our model to correctly identify the POIs. It is the probability of a POI being identified by our model. Below, we have the logistic regression precision and recall:

	precision	recall	f1-score	support
0.0	0.95	0.76	0.84	25
1.0	0.33	0.75	0.46	4
avg / total	0.86	0.76	0.79	29

0.1.3 References:

- [Advanced Machine Learning With Scikit Learn by Andreas Mueller](#)
- [Python Machine Learning by Sebastian Raschka](#)
- [Applied Predictive Modeling by Kuhn & Jhonson](#)