

Machine Learning Project

Enron Fraud Detectors using Enron Emails and Financial Data.

About Enron

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 public record, including tens of thousands of emails and detailed financial data for top executives. Most notable of which are Jefferey Skilling, Key Lay, and Fastow all have dumped large amounts of stock options, and they are all deemed guilty.

Goal of the Project

The goal of this project is to use financial and email data from Enron corpus - publicly made by US Federal Energy Regulatory Comission during its investgiation of Enron, which comprised email and financial data of 146 people (records), most of which are senior management of Enron, to come up with a predictive model that could spot an individual as a "Person of Interest" (POI).The "Person of Interest" here refers to a person who was involved in the fraud. The dataset contained 146 records with 1 labeled feature (POI), 14 financial features, 6 email feature. Within these record, 18 were labeled as a "Person Of Interest" (POI).

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. (All units are in US dollars)

Email Features: to_messages, email_address, from_poi_to_this_person, from_messages, from_this_person_to_poi, shared_receipt_with_poi. (units are generally number of emails messages; notable exception is 'email_address', which is a text string)

POI Label: poi. (boolean, represented as integer)

However, the dataset contains numerous people missing value for each feature which can be described as table below:

Feature	NaN per feature
Loan advances	142
Director fees	129
Restricted stock deferred	128
Deferred payment	107
Deferred income	97
Long term incentive	80

Bonus	64
Emails sent also to POI	60
Emails sent	60
Emails received	60
Emails from POI	60
Emails to POI	60
Other	53
Expenses	51
Salary	51
Excersised stock option	44
Restricted stock	36
Email address	35
Total payment	21
Total stock value	20

Outliers in the data

Through exploratory data analysis which included plotting of scatterplots using `matplotlib`, validating the list of keys in the data dictionary etc, I was able to figure out the following outliers -

`TOTAL` : Through visualising using scatter-plot matrix. We found `TOTAL` are the extreme outlier since it comprised every financial data in it.

`THE TRAVEL AGENCY IN THE PARK` : This must be a data-entry error that it didn't represent an individual.

Feature Scaling and Selection

Removing Features

By apply some intuition, the following facts could be formulated :

`email_address` can in no way differentiate between a POI and a non POI.

Features like `restricted_stock_deferred`, `director_fees`, `loan_advances` contains majority of their values as *NaN*.

Hence, these four features must be removed in order to gain better results.

Creating New Features

Instead of training the algorithm with the features like `from_this_person_to_poi` and `from_messages`, the ratio of number of messages from this person to poi and the total number of messages from this person would be more appropriate. Similar is the case with `from_poi_to_this_person` and `to_messages`.

Main purpose of composing ratio of POI message is we expect POI contact each other more often than non-POI and the relationship could be non-linear. The initial assumption behind these features is: the relationship between POI is much more stronger than between POI and non-POIs

New Features Created - `fraction_from_this_person_to_poi` and `fraction_from_poi_to_this_person`

Features Removed - `from_this_person_to_poi`, `from_poi_to_this_person`, `from_messages`, `to_messages`

Scaling Features

I have used `min_max_scaler` from preprocessing module of `scikit-learn` to scale the features. Scaling of features becomes crucial when the units of the features are not same.

Feature Selection

I used `scikit-learn SelectKBest` to select best 10 influential features and used those features for all the upcoming algorithm. Unsurprisingly, 8 out of 10 features are related to financial data and only 2 features are related to e-mail.

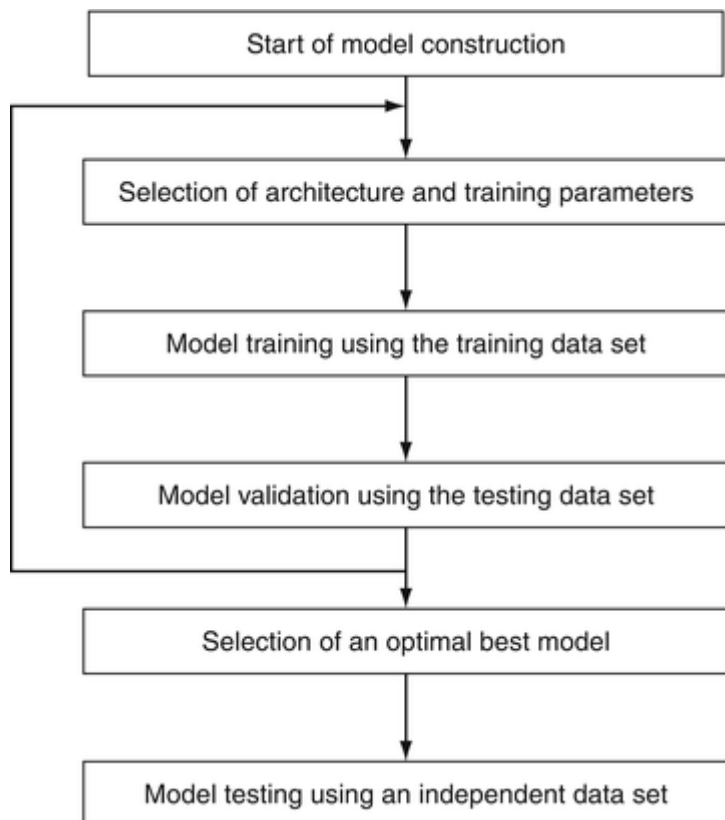
Selected Features	Score
<code>exercised_stock_options</code>	8.772
<code>total_stock_value</code>	11.458
<code>bonus</code>	24.815
<code>salary</code>	18.289
<code>deferred_income_</code>	9.922
<code>long_term_incentive</code>	16.409
<code>restricted_stock_</code>	20.792
<code>total_payments_</code>	0.225
<code>shared_receipt_with_poi_</code>	9.212
<code>fraction_from_this_person_to_poi_</code>	3.128

The six features in my POI identifier, which was a decision tree classifier. The first number is feature importance (from the decision tree classifier) and the second is feature score (from select k-best). The order of features is descending based on feature importance.

VALIDATION

In machine learning, model validation is referred to as the process where a trained model is evaluated with a testing data set. The testing data set is a separate portion of the same data set from which the training set is derived. The main purpose of using the testing data set is to test the generalization ability of a trained model .

Model validation is carried out after model training. As illustrated in Fig. 1, together with model training, model validation aims to find an optimal model with the best performance.



Cross-validation, sometimes called **rotation estimation**, is a model validation technique for assessing how the results of a statistical analysis will generalize to an independent data set. It is mainly used in settings where the goal is prediction, and one wants to estimate how accurately a predictive model will perform in practice. In a prediction problem, a model is usually given a dataset of *known data* on which training is run (*training dataset*), and a dataset of *unknown data* (or *first seen data*) against which the model is tested (*testing dataset*). The goal of cross validation is to define a dataset to "test" the model in the training phase (i.e., the *validation dataset*, in order to limit problems like overfitting¹ give an insight on how the model will generalize to an independent etc.

In my `poi_id.py` file, my data was separated into training and testing sets. The test size was 30% of the data, while the training set was 70%.

I also used `tester.py`'s cross validation as an alternate method to gauge my algorithm's performance. Because the Enron data set is so small, this type of cross validation was useful because it essentially created multiple datasets out of a single one to get more accurate results.

As described on [stats.stackexchange](#) ([link for k-nearest neighbours](#); [link for SVM](#)), normalization is required for:

k-nearest neighbours because the distance between the points drives the clustering that determines the nearest neighbours. If one feature has a much larger scale than another, the clustering will be driven by the larger scale and the distance between points on the smaller scale will be overshadowed. SVM because the distance between the separating plane and the support

vectors drives the algorithm's decision-making. If one feature has a much larger scale than another, it will dominate the other features distance-wise. Scaling isn't required for tree-based algorithms because the splitting of the data is based on a threshold value. This decision made based on this threshold value is not affected by different scales.

Selection Process

I used a univariate feature selection process, select k-best, in a pipeline with grid search to select the features. Select k-best removes all but the k highest scoring features. The number of features, 'k', was chosen through an exhaustive grid search driven by the 'f1' scoring estimator, intending to maximize precision and recall.

Splitting Data into Training and Testing Datasets

In order to split the data into training and testing datasets, I've used `train_test_split` from `cross_validation` module of `scikit-learn`.

```
features_train, features_test, labels_train, labels_test = \
    train_test_split(features, labels, test_size=0.3, random_state=42)
```

The parameter "test-size" adjusts the size of the testing dataset. In this case, 30% of the dataset is being used for testing.

Trying Different Classifier Algorithms and Evaluating Algorithms

The classifier I used in order to make the predictions are Support Vector Classifier, Decision Tree Classifier.

The decision tree classifier had a precision of 0.31336 and a recall of 0.59100, both above the 0.3 threshold. The SVM classifier (with features scaled) had a gaudy precision of 0.83333, but a poor recall of 0.06500. The k-nearest neighbors classifier had a precision of 0.32247 and a recall of 0.29200, both near but not both above the 0.3 threshold. As described by Udacity Coach Mitchell, the SVM classifier is not a good fit for the unbalanced classes in the Enron data set, i.e., the lack POIs vs. the abundance of non-POIs. The SVM overfit to the points in the training set, and poorly generalized to the test set resulting in a relatively small number of positive predictions of POIs.

Evaluation

The evaluation has been done by comparing the recall and precision obtained by using different classifier algorithms. `recall_score`, `precision_score` and `accuracy_score` has been used from `metrics` module of `scikit-learn`.

Parameter Tuning

Parameters tuning refers to the adjustment of the algorithm when training, in order to improve the fit on the test set. Parameter can influence the outcome of the learning process,

the more tuned the parameters, the more biased the algorithm will be to the training data & test harness. The strategy can be effective but it can also lead to more fragile models & overfit the test harness but don't perform well in practice

Tuning the parameters of an algorithm means adjusting the parameters in a certain way to achieve optimal algorithm performance. There are a variety of "certain ways" (e.g. a manual guess-and-check method or automatically with GridSearchCV) and "algorithm performance" can be measured in a variety of ways (e.g. accuracy, precision, or recall). If you don't tune the algorithm well, performance may suffer. The data won't be "learned" well and you won't be able to successfully make predictions on new data.

I used GridSearchCV for parameter tuning. As Katie Malone describes for Civis Analytics, GridSearchCV allows us to construct a grid of all the combinations of parameters, tries each combination, and then reports back the best combination/model.

For the chosen decision tree classifier for example, I tried out multiple different parameter values for each of the following parameters (with the optimal combination bolded). I used Stratified Shuffle Split cross validation to guard against bias introduced by the potential underrepresentation of classes (i.e. POIs).

criterion=['gini', 'entropy'] splitter=['best', 'random'] max_depth=[None, 1, 2, 3, 4]
min_samples_split=[1, 2, 3, 4, 25] class_weight=[None, 'balanced'] Note that a few of these chosen parameter values were different than their default values, which proves the importance of tuning.

Final Model of Choice

Precision is how often our class prediction (POI vs. non-POI) is right when we guess that class
Recall is how often we guess the class (POI vs. non-POI) when the class actually occurred In the context of our POI identifier, it is arguably more important to make sure we don't miss any POIs, so we don't care so much about precision. Imagine we are law enforcement using this algorithm as a screener to determine who to prosecute. When we guess POI, it is not the end of the world if they aren't a POI because we'll do due diligence. We won't put them in jail right away. For our case, we want high recall: when it is a POI, we need to make sure we catch them in our net.

For this assignment, I used precision & recall as 2 main evaluation metrics. The best performance belongs to Decision Tree Classifier (**precision: 0.314& recall: 0.587**) which is also the final model of choice. Precision refer to the ratio of true positive (predicted as POI) to the records that are actually POI while recall described ratio of true positives to people flagged as POI. Essentially speaking, with a precision score of 0.314, it tells us if this model predicts 100 POIs, there would be 41 people are actually POIs and the rest 49 are innocent. With recall score of 0.587, this model finds 45% of all real POIs in prediction. Due to the nature of the dataset, accuracy is not a good measurement as even if non-POI are all flagged, the accuracy score will yield that the model is a success.