

Enron Submission Free-Response Questions

By Lieby Cardoso
Udacity – Data Analyst Nanodegree

1. 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? [relevant rubric items: "data exploration", "outlier investigation"]

The energy company Enron enacted bankruptcy after executing several accounting and tax frauds. The purpose of this project is to use data analysis and machine learning techniques like features selection, preprocessing and evaluating models to predict if an employee has been involved with fraud. Those that were pointed out are called POIs (Person of Interest).

Dataset overview

- 146 data points, 21 features
- 18 (12%) POIs and 128 non-POI
- The original dataset consists of 1 POI identifier, 14 financial features, 6 emails related.
- Null values ratio:

| Feature | NaN ratio |
|---------------------------|-----------|
| loan_advances | 97% |
| restricted_stock_deferred | 88% |
| director_fees | 88% |
| deferral_paymen | 73% |
| deferred_income | 66% |
| long_term_incentive | 55% |
| bonus | 44% |
| to_messages | 41% |
| shared_receipt_with_poi | 41% |
| from_messages | 41% |
| from_this_person_to_poi | 41% |

| Feature | NaN ratio |
|-------------------------|-----------|
| from_poi_to_this_person | 41% |
| other | 36% |
| salary | 35% |
| expenses | 35% |
| exercised_stock_options | 30% |
| restricted_stock | 25% |
| total_payments | 14% |
| total_stock_value | 14% |
| poi | 0 |
| email_address | 0 |

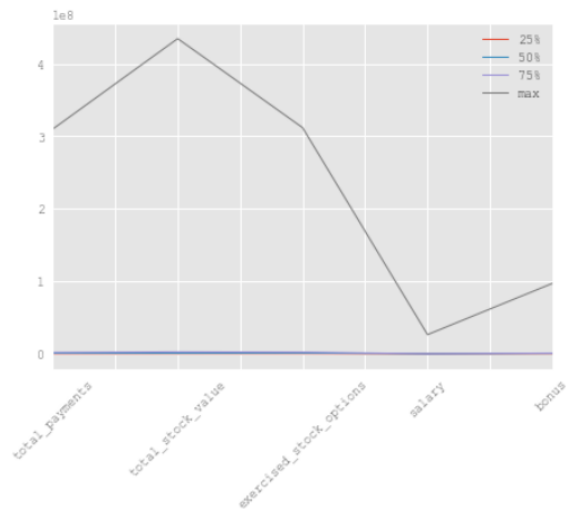
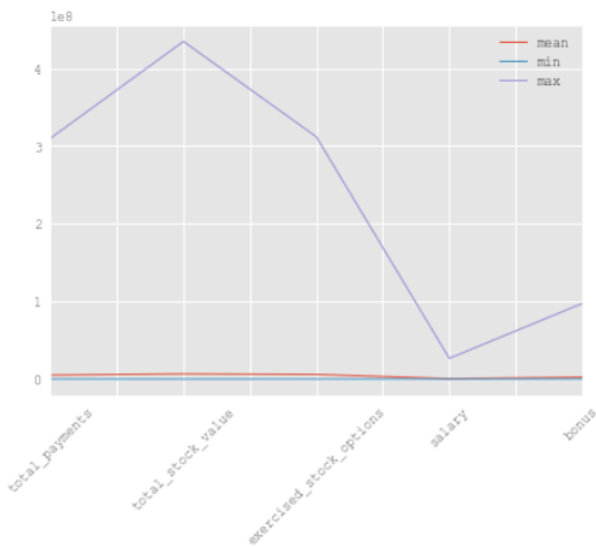
Outliers Investigation

The data set used is composed of financial information, the emails sent and received by the people of interest (POI) and the identifier of whether the key person is a POI.

Enron Submission Free-Response Questions

By Lieby Cardoso

Udacity – Data Analyst Nanodegree



According to Forbes, "the stated goal of its board of directors was to pay executives in the 75th percentile of its peer group". And we could prove this when we searched for employees with total payments above the 3rd quartile. In the graph above we can see that the max values for the selected features is well above the 3rd quartile and the mean. Jeffrey Skilling (CEO) and Kenneth Lay (chairman and CEO) have many above average values, but justified by the importance of their positions.

When I was looking for outliers I found a POI called TOTAL. It was an outlier whose content is the sum of the other values, I deleted this record. The second largest value receive was Kenneth Lay, but, as you know this make perfect sense since he was the CEO and chairman of Enron.

Other records analyzed and deleted:

- "THE TRAVEL AGENCY IN THE PARK" has been deleted. Although the agency had received \$ 350,000 in payments 2 days before Enron's bankruptcy and Sharon Lay (sister of Kenneth Lay) owned 50% of the company, this record will not be considered a POI.
- "LOCKHART EUGENE E" doesn't have any value assign.
- "CHAN RONNIE" had stock and income put off to a later time causing a total payment equal a 0, and he had no message sent or received. Since he is not a POI, his presence on the data set was not justified.

Records analyzed and kept:

- "POWERS WILLIAM" had stock or income put off to a later time causing a total payment equal a 0, but has values assigned to shared_receipt_with_poi.

Enron Submission Free-Response Questions

By Lieby Cardoso
Udacity – Data Analyst Nanodegree

- "HORTON STANLEY C" CEO of Enron Transportation Services had no salary and bonus amount associated with it, but had deferred payments and stock assigned.
- "BHATNAGAR SANJAY" and "BELFER ROBERT" had wrong values assign as total payments and stock, So I fixed them up.

To see how I clean the data go to:

https://github.com/liebycardoso/ML_Enron/blob/master/Enron_Data_Analysis.html

2. 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. (You do not necessarily have to use it in the final analysis, only engineer and test it.) In your feature selection step, if you used an algorithm like a decision tree, please also give the feature importance's of the features that you use, and if you used an automated feature selection function like SelectKBest, please report the feature scores and reasons for your choice of parameter values. [relevant rubric items: "create new features", "properly scale features", "intelligently select feature"]

After clearing the data, I created two new features:

1. message_poi_ratio: ratio of emails related to the POI, both sent and received;
2. message_others_ratio: ratio of messages not related with a POI.

Creating new variables made e-mail features unnecessary. I did not use them to avoid duplicate data. Dropped features: "email_address", "from_messages", "from_poi_to_this_person", "from_this_person_to_poi", "shared_receipt_with_poi", "to_messages". In the same way, "total_payments" and "total_stock_value" are totalizers and were excluded because their data are distributed in the financial variables.

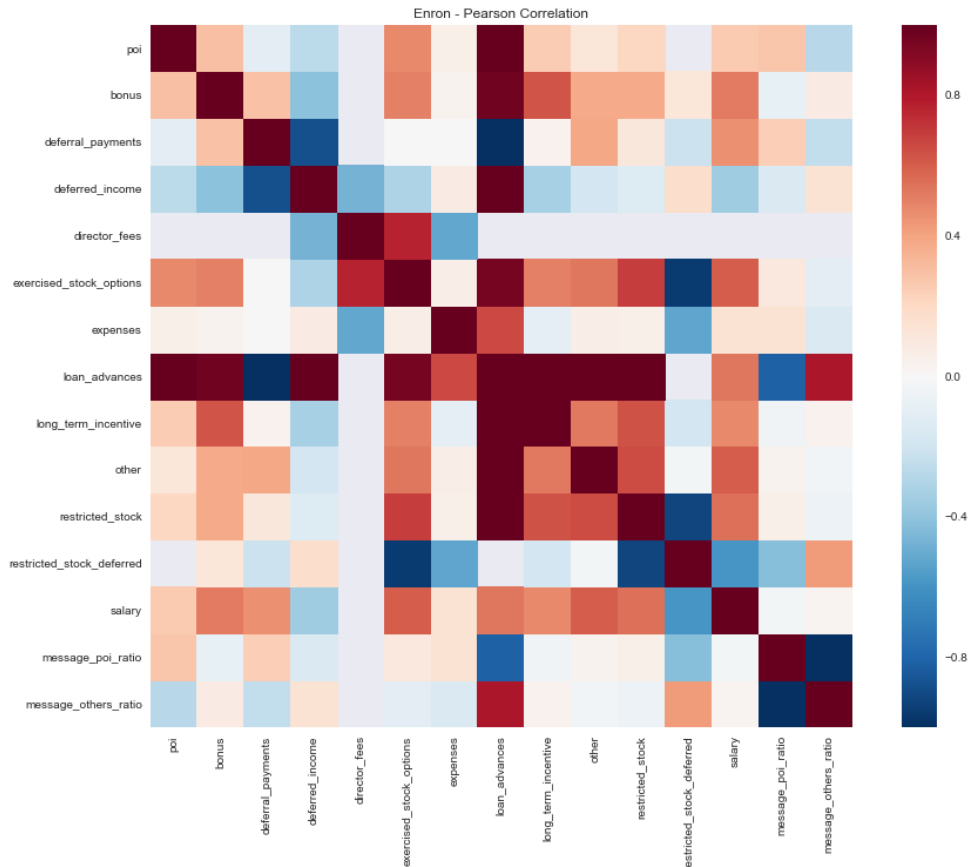
Accuracy, precision and Recall scores for Adaboost classifier with original dataset and after I created two features.

| | Accuracy | Precision | Recall |
|---|----------|-----------|---------|
| <i>Performance original clean dataset</i> | 0.71340 | 0.27873 | 0.72400 |
| <i>Performance new features</i> | 0.80280 | 0.36377 | 0.63950 |

Once the dataset was ready, the first thing I did was apply the pandas.DataFrame's corr() function on all variables to get the correlation between them.

Enron Submission Free-Response Questions

By Lieby Cardoso
Udacity – Data Analyst Nanodegree



The feature POI had the strongest relationship with the variables Loan_advances (0.99), exercised_stock_options (0.48), bonus (0.30), message_poi_ratio (0.27), message_others_ratio (-0.27), salary (0.26) and deferred_income (0.26)

I used the SelectKBest class to score features using function Anova and Chi2. I created a comparison between the Kbest scores on the old dataset and the dataset with the new features. As you can see the feature message_poi_ratio was among the top 10 scores.

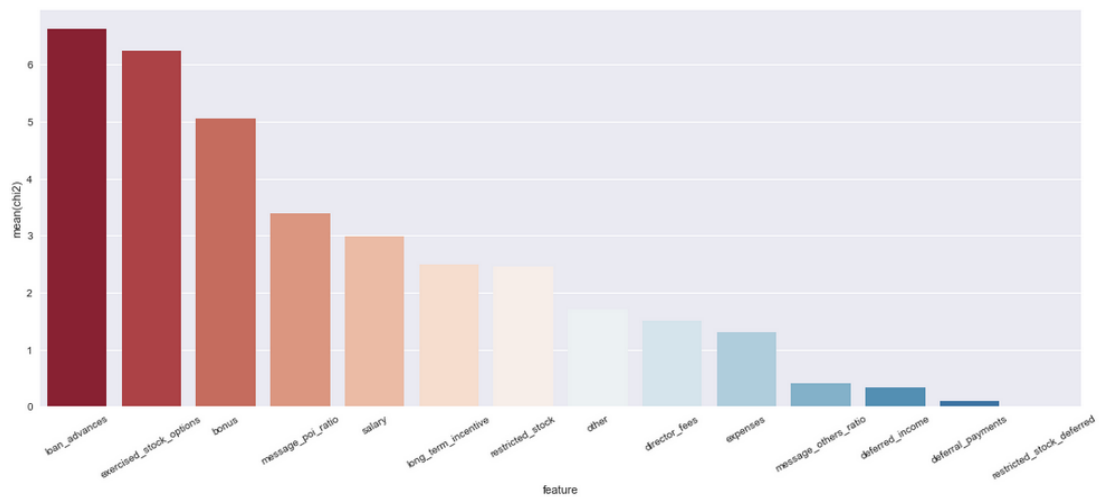
| feature | Anova | Chi2 | Anova Original Dataset |
|--------------------------------|-----------|----------|------------------------|
| <i>exercised_stock_options</i> | 22.087532 | 6.246159 | 22.087532 |
| <i>bonus</i> | 20.524645 | 5.048256 | 20.524645 |
| <i>salary</i> | 18.003740 | 2.989183 | 18.003740 |
| <i>deferred_income</i> | 11.320185 | 0.338413 | 11.320185 |
| <i>message_poi_ratio</i> | 9.816852 | 3.371674 | Absent |
| <i>long_term_incentive</i> | 9.772104 | 2.497366 | 9.772104 |
| <i>restricted_stock</i> | 8.694888 | 2.463442 | 8.694888 |
| <i>loan_advances</i> | 7.125382 | 6.634816 | 7.125382 |
| <i>expenses</i> | 5.287549 | 1.293600 | 5.287549 |

Enron Submission Free-Response Questions

By Lieby Cardoso

Udacity – Data Analyst Nanodegree

| | | | |
|----------------------------------|----------|----------|----------|
| <i>other</i> | 4.143788 | 1.703679 | 4.143788 |
| <i>director_fees</i> | 1.972788 | 1.508957 | 1.972788 |
| <i>message_others_ratio</i> | 1.473822 | 0.392669 | Absent |
| <i>restricted_stock_deferred</i> | 0.767702 | 0.008756 | 0.767702 |
| <i>deferral_payments</i> | 0.236711 | 0.094344 | 0.236711 |
| <i>shared_receipt_with_poi</i> | Absent | Absent | 8.432635 |
| <i>from_poi_to_this_person</i> | Absent | Absent | 5.446687 |
| <i>from_this_person_to_poi</i> | Absent | Absent | 2.338836 |
| <i>to_messages</i> | Absent | Absent | 1.594256 |
| <i>from_messages</i> | Absent | Absent | 0.175383 |

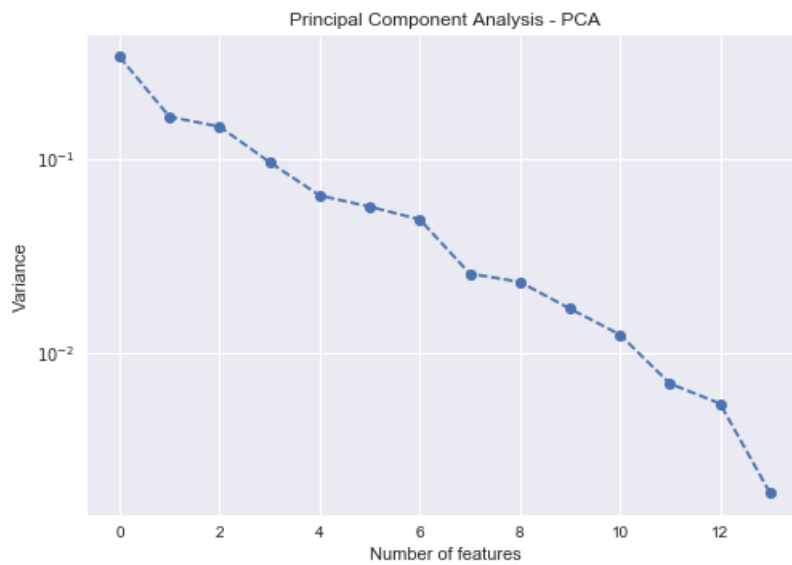


We have a very similar result between the correlation of Pearson and SelectKbest with Chi-squared stats on a scaling (0.1) data.

Another approach was to reduce the dimensionality of the dataset, I used the Principal Component Analysis (PCA) to identify and keep the most influential dimensions.

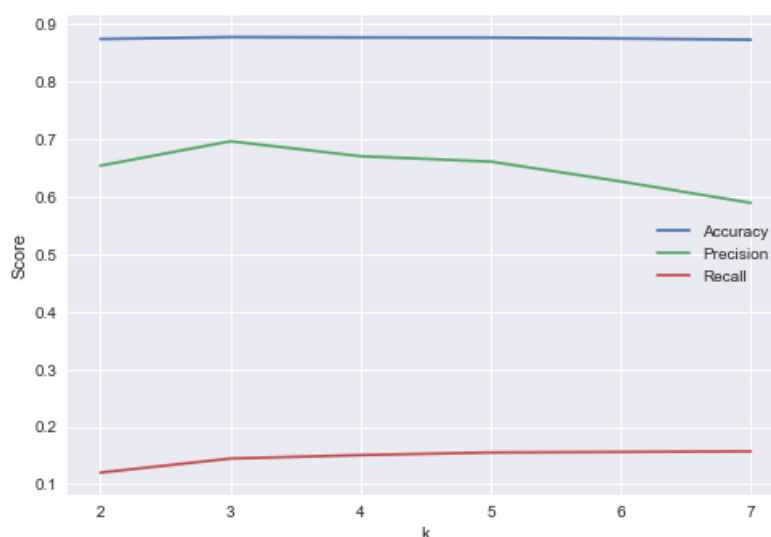
Enron Submission Free-Response Questions

By Lieby Cardoso
Udacity – Data Analyst Nanodegree



After six components, we have a drop on variance ratio.

During the feature selection step, I evaluated the pearson correlation, the PCA variance ratio and the K scores. I concluded that the ideal value of K or n_components would be something around 5 and 6. For the LogisticRegression and KNeighborsClassifier classifiers the selection of features with Kbest ($k = 5$) and PCA ($n = 5$) was fundamental, but unfortunately for the final chosen classifier I could not reach a precision greater than 0.30. In the graph below you can see the variation of the scores for each k value on Adaboost:



Enron Submission Free-Response Questions

By Lieby Cardoso
Udacity – Data Analyst Nanodegree

3. What algorithm did you end up using? What other one(s) did you try? How did model performance differ between algorithms? [relevant rubric item: "pick an algorithm"]

The answer I was pursuing to know, based on my data, was if a given person is a POI or not. Since there are only two possible answers, we can say that this is a classification task and to help me define the way forward, I have tested multiple classifiers and their default parameters.

| Model | Precision | Recall | Accuracy |
|----------------------------|-----------|--------|----------|
| KNeighbors | 0.581633 | 0.0570 | 0.868800 |
| RandomForest | 0.460692 | 0.1465 | 0.863333 |
| ExtraTreesClassifier | 0.423295 | 0.1490 | 0.859467 |
| AdaBoost | 0.418319 | 0.3060 | 0.850733 |
| NearestCentroid | 0.376596 | 0.2655 | 0.843467 |
| GradientBoostingClassifier | 0.263557 | 0.2260 | 0.812600 |
| DecisionTree | 0.257143 | 0.2475 | 0.804333 |
| LogisticRegression | 0.164218 | 0.1900 | 0.763067 |
| Naive Bayes | 0.219135 | 0.6230 | 0.653733 |

Based on the accuracy, precision and recall value, I chose AdaBoost, GradientBoostingClassifier and NearestCentroid to improve performance. My goal was to achieve a precision and recall value higher than 0.3. In this scenario, with AdaBoost I achieve my goals even with default parameters.

Scores obtained after tuning:

| Model | Precision | Recall | Accuracy |
|----------------------------|-----------|---------|----------|
| KNeighbors | 0.39456 | 0.37700 | 0.83980 |
| RandomForest | 0.43103 | 0.13750 | 0.86080 |
| ExtraTreesClassifier | 0.22082 | 0.68300 | 0.63640 |
| AdaBoost | 0.36377 | 0.63950 | 0.80280 |
| NearestCentroid | 0.35857 | 0.52800 | 0.81113 |
| GradientBoostingClassifier | 0.31915 | 0.23250 | 0.83153 |
| DecisionTree | 0.26790 | 0.75200 | 0.69293 |
| LogisticRegression | 0.30177 | 0.57200 | 0.76647 |

At some point in the Project I was curious about how much I could improve other classifiers, so I continued the tuning process. Only Gaussian Naive Bayes was left out, because this estimator does not have parameters to work with, just priors.

Enron Submission Free-Response Questions

By Lieby Cardoso
Udacity – Data Analyst Nanodegree

I got three models that stood out, both KNeighbors and NearestCentroid had a superior performance when compared to others, but the chosen one was AdaBoost which had precision values like KNeighbors and NearestCentroid, but a higher recall value.

To see how I've decide for this classifier, check:

https://github.com/liebycardoso/ML_Enron/blob/master/Enron_Model_Selection.html

4. 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? (Some algorithms do not have parameters that you need to tune -- if this is the case for the one you picked, identify and briefly explain how you would have done it for the model that was not your final choice or a different model that does utilize parameter tuning, e.g. a decision tree classifier). [relevant rubric item: "tune the algorithm"]

With the gridsearchcv function an exhaustive search of the best parameter for the estimator is done, sometimes in this search for the best model you set a value that increase the Recall metric and decrease Precision. To be successful in this task the range of parameters must be carefully selected, in my test I had cases where my selection of parameters was incorrect and the tuning model returned metrics worse than its version with default parameters.

I tried reducing the dimensionality of the data with PCA and selecting the features with SelectkBest, but in some cases, there was no improvement as I would like. I used MinMaxScale to transforms features by scaling then to a range between 0 and 1.

First, I tune the decisionTreeClassifier with the following parameters in order to fine-tune its performance with Adaboost.

| Classifier: DecisionTreeClassifier | | |
|------------------------------------|---------------------|----------|
| Parameter | Options | Picked |
| max_features | [2, 6, 10] | 2 |
| min_samples_split | [0.1, .50, .90] | 0.5 |
| splitter | ['best', 'random'] | random |
| min_samples_leaf | [1, 3, 10,15] | 1 |
| class_weight | ['balanced', None] | balanced |
| criterion | ["gini", "entropy"] | gini |

Enron Submission Free-Response Questions

By Lieby Cardoso

Udacity – Data Analyst Nanodegree

In the sequence, I used GridSearchCV to tune the Adaboost with the parameters listed below.

| Classifier: Adaboost | | |
|----------------------|------------------------------|--------|
| Parameter | Options | Picked |
| n_estimators | [150,200] | 150 |
| learning_rate | [0.01, 0.05, 0.1, 1.0] | 0.1 |
| algorithm | ['SAMME.R', 'SAMME'] | SAME |
| base_estimator | DecisionTreeClassifier tuned | |

The Adaboost achieved optimal result with scaled features on a range of 0 and 1, algorithm SAME, a learning rate of 0.1 e number of estimators equal 150.

Final pipeline:

```
Pipeline(steps=[
    ('minmaxer', MinMaxScaler(copy=True, feature_range=(0, 1))),
    ('classifier', AdaBoostClassifier(algorithm='SAMME',
    base_estimator=DecisionTreeClassifier(class_weight='balanced',
    criterion='gini', max_depth=None, max_features=2, max_leaf_nodes=None,
    random_state=42, splitter='random'),
    learning_rate=0.1,
    n_estimators=150,
    random_state=42))])
```

5. What is validation, and what's a classic mistake you can make if you do it wrong? How did you validate your analysis? [relevant rubric item: "validation strategy"]

To ensure that our data model can predict results, even in an unknown dataset, we need to separate our data sample into a training set and a test set. Not doing this division can lead us to an overfitting situation, which is when the model foresees perfectly what has already seen, but fails to predict what has not been tested.

While I was testing the models, I split the dataset using `cross_validation.train_test_split` and set the test set size in 30%. In a small dataset like ours, where we do not have many examples for learning, it is important to ensure that by separating the data samples into training and testing we have in each sample a representation of the classes that are important to us. One way to do this is to use the `StratifiedShuffleSplit` method that randomizes the order of the records before doing the split and ensures that each dataset contains a similar number of classes (examples) in each sample. I used this validator in the final test (Script `tester.py`) with `n_splits=1000` so that the process of fitting and training was repeated 1000 times in a way I could achieve the most accurate performance.

Enron Submission Free-Response Questions

By Lieby Cardoso
Udacity – Data Analyst Nanodegree

6. 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. [relevant rubric item: "usage of evaluation metrics"]

Once a model is considered ready it is important to choose the statistical metrics that can help us estimate how accurately and precise our model is. I started this project using the F1 metric as a reference, but it is a balance between precision and recall and it ended up returning a high value when one of the measurements was high and I needed both recall and precision to be over 0.3. This project focused on precision, recall and accuracy.

The model chosen after tuning was AdaBoostClassifier that returned the most balanced values of precision, recall and accuracy.

Achieved values:

| Classifier | Precision | Recall | Accuracy |
|---|-----------|---------|----------|
| AdaBoostClassifier before tuning | 0.418319 | 0.3060 | 0.850733 |
| AdaBoostClassifier after tuning | 0.36377 | 0.63950 | 0.80280 |

As you may realize after tuning the model has lost a bit of accuracy and precision, but has achieved a significant increase in recall value. Which is quite significant for me because it tells me that the model has greatly improved its ability to properly point a POI.

The performance metrics of our model tells us that it is capable of:

- For POI records predicted, 36% of them are POIs (Precision)
- Among all available records it can identify 64% of POIs (Recall)
- Correctly predict 80% of the values (Accuracy)

Identifying whether a particular employee is related to a fraud scheme is a very important task and that if misclassified can have profound consequences for the people listed. For this type of project would be very important to have a dataset with more than 18 POIs for training.