

Udacity Project Machine Learning - Enron Mail Dataset

Introduction

The goal of the project is to study a dataset of emails of a company called Enron. This company went bankrupt after a the fraud was detected.

"Enron's complex financial statements were confusing to shareholders and analysts. In addition, its complex business model and unethical practices required that the company use accounting limitations to misrepresent earnings and modify the balance sheet to indicate favorable performance." source:

https://en.wikipedia.org/wiki/Enron_scandal (https://en.wikipedia.org/wiki/Enron_scandal)

The information if a person is POI or not is provided by Udacity.

Project goal

The goals of the project are:

- Get familiar with dataset, clean dataset, check dataset for outliers
- Visualize dependencies of features
- Create new features (feature engineering)
- Set up ML model to indentify if a person was a "POI" or not
- Select best features to achieve a good presicion and recall score
- Tune selected model
- Validate model against test data

Validation strategy

The validation of a model is important, because we need to know how good the model is at identifying POI from the dataset. We will use a classical train/test split of the entire dataset, which divides the data set into 70% data used to train the ML model and 30% of testing/validation data.

This will allow us to check the ML model performance on data, which the model has never seen before. If we would train the model on the complete dataset, we will in most cases get a perfect validation, which is misleading. That's why validation of a ML model needs to be done on a independent set of data.

The validity of the classifier model is measured by the precision and recall score. In this way we can justify the feature selection and the model hyper-parameter tuning.

"Precision is the number of correct positive classifications divided by the total number of positive labels assigned. In other words, it is the fraction of persons of interest predicted by the algorithm that are truly persons of interest. Mathematically precision is defined as" (Ref. 1)

- $\text{precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$

"Recall is the number of correct positive classifications divided by the number of positive instances that should have been identified. In other words, it is the fraction of the total number of persons of interest in the data that the classifier identifies. Mathematically, recall is defined as" (Ref. 1)

- $\text{recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$

https://en.wikipedia.org/wiki/Precision_and_recall (https://en.wikipedia.org/wiki/Precision_and_recall)

Script overview

The script "poi_id.py" and the created *.pkl files are located at the sub-folder "final_project". The project evaluator will test these using the tester.py script.

```
In [172]: 1 #!/usr/bin/python
2
3 import sys
4 import pickle
5 sys.path.append("../tools/")
6
7 import matplotlib.pyplot as plt
8 import numpy as np
9 import pandas as pd
10 import seaborn as sns
11
12 from sklearn import preprocessing
13 from sklearn.preprocessing import MinMaxScaler
14 from sklearn.cross_validation import train_test_split
15 from time import time
16 from sklearn.naive_bayes import GaussianNB
17 from sklearn.metrics import accuracy_score, precision_score, recall_score,
18 from sklearn.grid_search import GridSearchCV
19
20 from feature_format import featureFormat, targetFeatureSplit
```

Task 1: Feature Selection

```
In [173]: 1 ### Load the dictionary containing the dataset
2 with open("final_project_dataset.pkl", "r") as data_file:
```

```
In [174]: 1 ### Task 1: Select what features you'll use.
2 ### features_list is a list of strings, each of which is a feature name.
3 ### The first feature must be "poi".
4 features_list_init = ["poi", "salary", "bonus", 'from_poi_to_this_person', 'to_poi_to_this_person', 'total_payments', 'loan_advances', 'restricted_stock_deferred', 'total_stock_value', 'exercised_stock_options', 'expenses', 'exercised_stock_options', 'other', 'from_this_person_to_poi', 'poi_email_address', 'director_fees']
5
6 features_list = ["poi", "salary", "bonus", 'from_poi_ratio', 'to_poi_ratio', 'deferral_payments', 'total_payments', 'loan_advances', '#', 'deferred_income', 'total_stock_value', 'expenses', 'exercised_stock_options', 'other', 'long_term_incentive', 'shared_rec_ratio', 'restricted_stock_deferred', '#', 'director_fees']
7
8 #number of features used:
9 print("number of features:", len(features_list))
10
11 ('number of features:', 15)
```

Data Exploration

```
In [175]: 1 # print names of all 146 individuals:
```

```
In [176]: 1 # print entries of the first person:
2 print(data_dict['PHILLIP ALLEN RUTLEDGE'])
3
4 {'salary': 201955, 'to_messages': 2902, 'deferral_payments': 2869717, 'total_payments': 4484442, 'exercised_stock_options': 1729541, 'bonus': 4175000, 'restricted_stock': 126027, 'shared_receipt_with_poi': 1407, 'restricted_stock_deferred': -126027, 'total_stock_value': 1729541, 'expenses': 13868, 'loan_advances': 'NaN', 'from_messages': 2195, 'other': 152, 'from_this_person_to_poi': 65, 'poi_email_address': False, 'director_fees': 'NaN', 'deferred_income': -3081055, 'long_term_incentive': 304805, 'email_address': 'phillip.allen@enron.com', 'from_poi_to_this_person': 47}
```

Data Exploration

First, I will import the dict into a pandas Dataframe, since it will make the data exploration and clean up much easier for me.

According to the documentation of the enron mail dataset the NAN values of financial data are related to a 0. This is not true for the email address, but replacing a NAN with a 0 here will not have an influence on results, since the email address is not a candidate for a feature.

In [177]:

Number of persons within the dataset: 146

146, but 1 value is the "total" row, which is removed later.

In [178]:

```
1 df = pd.DataFrame(data_dict)
2 df = df.T
```

Out[178]:

	bonus	deferral_payments	deferred_income	director_fees	email_address	exercised_stock_options
ALLEN PHILLIP K	4175000	2869717	-3081055	NaN	phillip.allen@enron.com	
BADUM JAMES P	NaN	178980	NaN	NaN		NaN
BANNANTINE JAMES M	NaN	NaN	-5104	NaN	james.bannantine@enron.com	
BAXTER JOHN C	1200000	1295738	-1386055	NaN		NaN
BAY FRANKLIN R	400000	260455	-201641	NaN	frank.bay@enron.com	

5 rows × 7 columns

I replace the string "NaN" with np.nan and count the nan per feature. Some features contain many nan, but as written above this actually means 0. For this reason I replace np.nan with 0.

In [179]:

```
1 df = df.replace('NaN', np.nan)
```

Out[179]:

```
bonus          64
deferral_payments  107
deferred_income   97
director_fees   129
email_address    35
exercised_stock_options  44
expenses        51
from_messages    60
from_poi_to_this_person  60
from_this_person_to_poi  60
loan_advances   142
long_term_incentive  80
other          53
poi             0
restricted_stock  36
restricted_stock_deferred  128
salary         51
shared_receipt_with_poi  60
to_messages     60
total_payments  21
total_stock_value  20
dtype: int64
```

In [180]:

In [181]:

Out[181]:

	bonus	deferral_payments	deferred_income	director_fees	email_address	exercised_stock_options
ALLEN PHILLIP K	4175000.0	2869717.0	-3081055.0	0.0	phillip.allen@enron.com	
BADUM JAMES P	0.0	178980.0	0.0	0.0		0
BANNANTINE JAMES M	0.0	0.0	-5104.0	0.0	james.bannantine@enron.com	
BAXTER JOHN C	1200000.0	1295738.0	-1386055.0	0.0		0
BAY FRANKLIN R	400000.0	260455.0	-201641.0	0.0	frank.bay@enron.com	

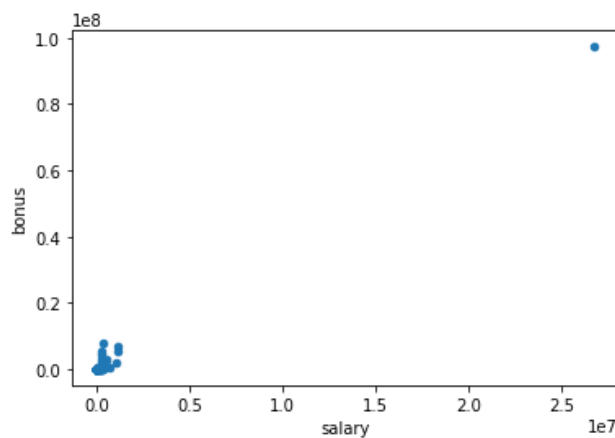
5 rows × 21 columns

Task 2: Remove outliers

The total entry is obviously an outlier, which will be dropped from the dataframe

In [182]:

Out[182]: <matplotlib.axes._subplots.AxesSubplot at 0x7f5e45ac9710>



In [183]:

Out[183]:

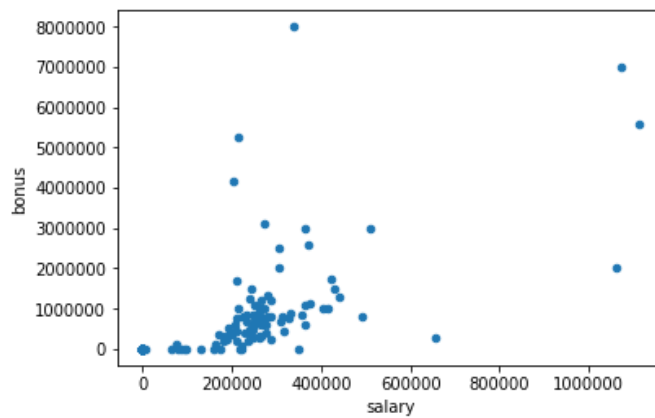
	bonus	deferral_payments	deferred_income	director_fees	email_address	exercised_stock_options
TOTAL	97343619.0	32083396.0	-27992891.0	1398517.0		0

1 rows × 21 columns

The "total" row is an obvious outlier and will be removed.

```
In [184]: 1
          2 df = df.drop(['TOTAL'])
          3 data_dict.pop('TOTAL', 0)
```

Out[184]: <matplotlib.axes._subplots.AxesSubplot at 0x7f5e4589c650>

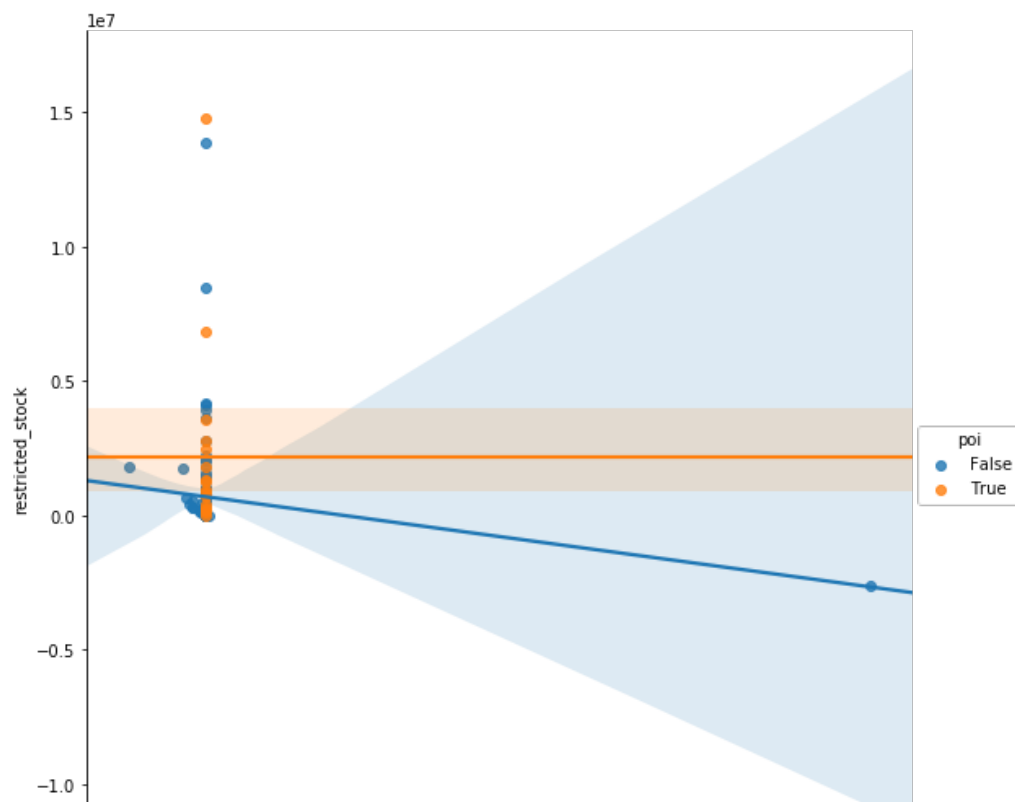
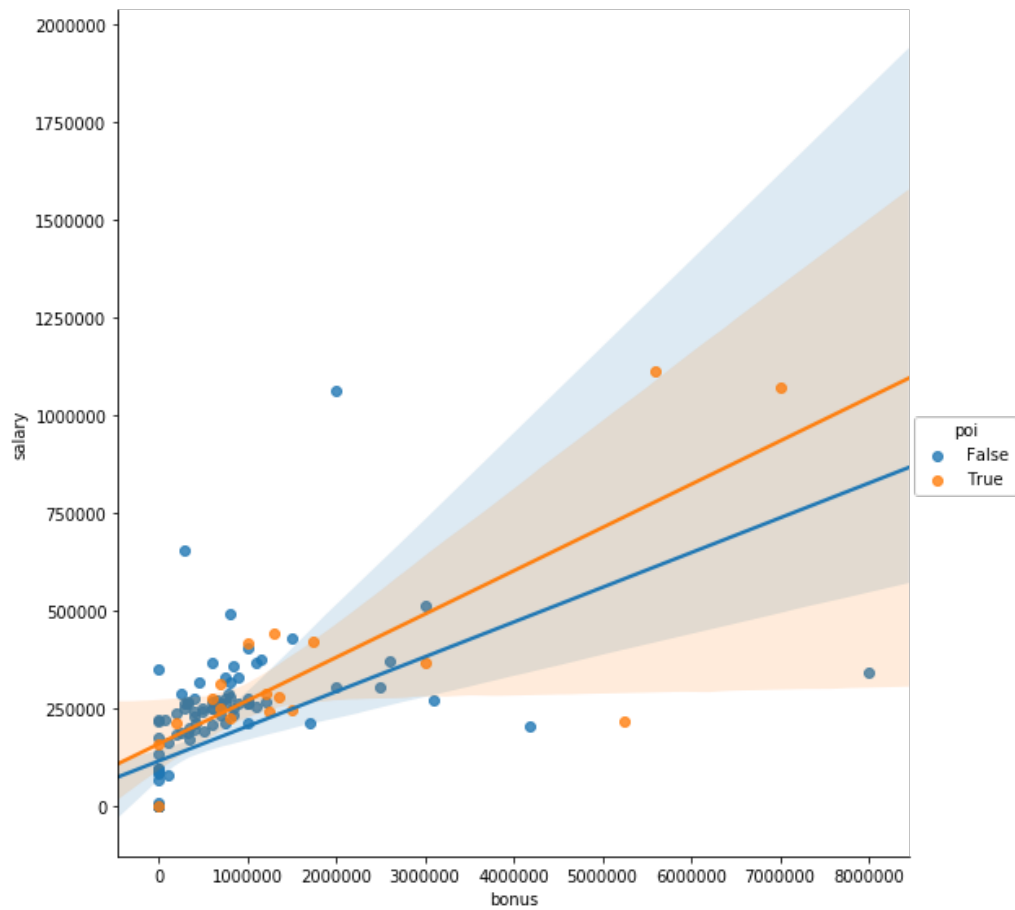


```
In [185]: 1 # count number of POI/non-POI
```

Out[185]: False 127
 True 18
 Name: poi, dtype: int64

```
In [186]: 1 #sns.pairplot(df, vars=['salary', 'total_payments'], hue= 'poi', size = 5)
          2 sns.lmplot(data = df, x = 'bonus', y = 'salary', hue = 'poi', size = 8)
```

Out[186]: <seaborn.axisgrid.FacetGrid at 0x7f5e4591ad10>



```
In [187]: 1 # looking at the graph above, I think the feature "restricted_stock_deferred"
2 # Indeed, by removing this feature I increased the score from 0.66 to 0.91.
3 df[df['restricted_stock_deferred']==df['restricted_stock_deferred'].max()]
```

```
Out[187]:
```

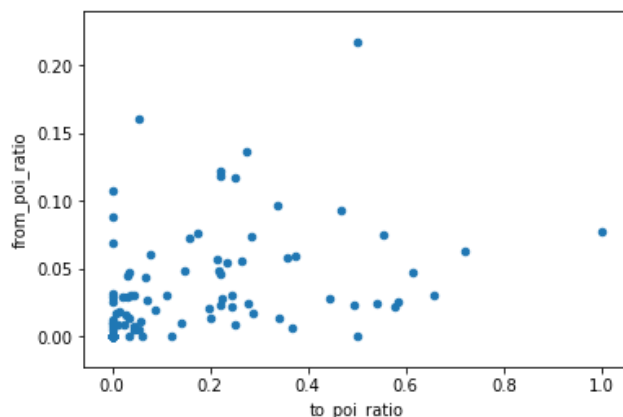
	bonus	deferral_payments	deferred_income	director_fees	email_address	exercised
BHATNAGAR SANJAY	0.0	0.0	0.0	137864.0	sanjay.bhatnagar@enron.com	

1 rows × 7 columns

Task 3: Create new feature(s)

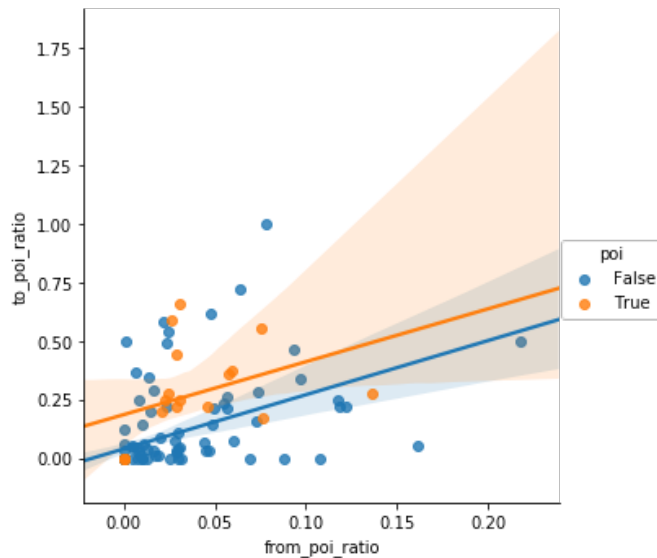
```
In [188]: 1 ### Task 3: Create new feature(s)
2
3 #Note that I only create some new features for the sake of the project submission
4 #At the end I will use a decision tree classifier, so scaling does not influence
5
6 df['from_poi_ratio'] = df['from_poi_to_this_person'] / df['to_messages']
7 df['to_poi_ratio'] = df['from_this_person_to_poi'] / df['from_messages']
8 df['shared_rec_ratio'] = df['shared_receipt_with_poi'] / df['to_messages']
9 df = df.fillna(0)
10 df.plot('to_poi_ratio', 'from_poi_ratio', kind = 'scatter')
11
```

```
Out[188]: <matplotlib.axes._subplots.AxesSubplot at 0x7f5e456eab10>
```



In [189]:

Out[189]: <seaborn.axisgrid.FacetGrid at 0x7f5e458b0d10>



In the following lines of code I convert the pandas dataframe back to a dict, since the tester.py script expects this format.

In [190]:

```

1  ### Store to my_dataset for easy export below.
2  #my_dataset = data_dict
3  my_dataset = df.to_dict(orient='index')
4
5  ### Extract features and labels from dataset for local testing
6  #data = featureFormat(my_dataset, features_list_init, sort_keys = True)
7  #labels, features = targetFeatureSplit(data)
8  labels_df = df['poi']
9  features_df = df[features_list].drop(['poi'], axis = 1)

```

Task 4: Try a variety of classifiers


```

In [191]: 1  ### Task 4: Try a variety of classifiers
2  ### Please name your classifier clf for easy export below.
3
4  features_train, features_test, labels_train, labels_test = \
5      train_test_split(features_df, labels_df, test_size=0.3, random_state=42)
6  #features_train, features_test, labels_train, labels_test = \
7  #      train_test_split(features, labels, test_size=0.3, random_state=42)
8
9
10 from sklearn.svm import SVC
11 clf = SVC()
12 clf.fit(features_train, labels_train)
13 pred = clf.predict(features_test)
14 accuracy = accuracy_score(labels_test, pred)
15 prec = precision_score(labels_test, pred)
16 recall = recall_score(labels_test, pred)
17 print "SVC accuracy score:", "%.2f" % round(accuracy,3) , "precision:", "%.2f"
18
19 from sklearn.tree import DecisionTreeClassifier
20 clf = DecisionTreeClassifier(random_state=2, max_features=5)
21 clf_select = clf
22 clf.fit(features_train, labels_train)
23 clf_select.fit(features_train, labels_train)
24 score = clf.score(features_test, labels_test)
25 pred = clf.predict(features_test)
26 accuracy = accuracy_score(labels_test, pred)
27 prec = precision_score(labels_test, pred)
28 recall = recall_score(labels_test, pred)
29 print "DTC accuracy score:", "%.2f" % round(accuracy,3) , "precision:", "%.2f"
30
31 from sklearn.ensemble import RandomForestClassifier
32 clf = RandomForestClassifier(random_state=2)
33 clf.fit(features_train, labels_train)
34 score = clf.score(features_test, labels_test)
35 pred = clf.predict(features_test)
36 accuracy = accuracy_score(labels_test, pred)
37 prec = precision_score(labels_test, pred)
38 recall = recall_score(labels_test, pred)
39
40 print "RFC accuracy score:", "%.2f" % round(accuracy,3) , "precision:", "%.2f"
41
42 print confusion_matrix(labels_test, pred)

```

SVC accuracy score: 0.91 precision: 0.00 recall: 0.00

DTC accuracy score: 0.89 precision: 0.40 recall: 0.50

RFC accuracy score: 0.91 precision: 0.50 recall: 0.50

```
[[38  2]
 [ 2  2]]
```

/home/niklas/Documents/anaconda3/envs/python2/lib/python2.7/site-packages/sklearn/metrics/classification.py:1135: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no predicted samples.
'precision', 'predicted', average, warn_for)

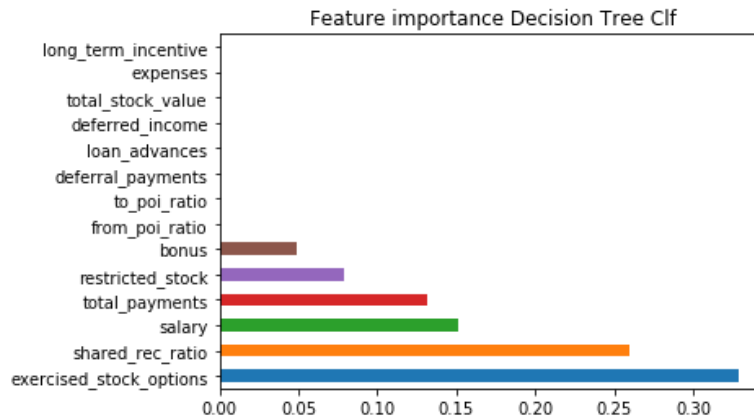
pick an algorithm

By manually changing the hyper-parameter 'max_features' to 8 I could increase the precision score from 0.25 to 0.31 and the recall score from 0.25 to 0.33.

Initially, the DTC did perform best, so I have chosen this model for further fine tuning. Note that the initial precision and recall score of the RFC was worse than the DTC. This changed after removing the 'directors fee' feature from the list.

```
In [192]: 1 #plot the feature importances of random decision tree classifier
          2 (pd.Series(clf_select.feature_importances_, index=features_train.columns)
          3 .nlargest(15))
```

Out[192]: <matplotlib.axes._subplots.AxesSubplot at 0x7f5e45eb5990>



Task 5: Tune the classifier

The precision and recall rate of the decision tree classifier is quite promising, and I will try to tune it in this chapter.

```
In [193]: 1 from sklearn.pipeline import Pipeline
          2 from sklearn.grid_search import GridSearchCV
          3 from sklearn.feature_selection import SelectKBest, f_classif
```

Automated Feature Selection with GridSearchCV with SelectKBest

I will now use GridSearchCV with SelectKBest to search for the best features for the decision tree classifier. As proposed by the reviewer, I combine the selection of features and the algorithm using a pipeline. In this way the best features are selected in an automated way. The GridSearchCV tunes the "number of features to be selected" and the hyperparameter of the estimator, by selecting the parameters that give the best score on validation data.

```
In [194]: 1 n_features = np.arange(1, len(features_list))
2 kbest = SelectKBest(f_classif)
3 #param_grid = [{'select_features_k': n_features}]
4 # Use GridSearchCV to automate the process of finding the optimal number of
5 #tree_clf= GridSearchCV(pipe, param_grid=param_grid, scoring='f1', cv = 10)
6 #tree_clf.fit(features_train,labels_train)
7
8 pipeline = Pipeline([('kbest', kbest), ('classify', DecisionTreeClassifier())])
9 grid_search = GridSearchCV(pipeline, {'kbest__k': [6,8,10,12,14], 'classify__min_samples_split': [2,4,6,8,10,12,14]}, scoring='f1')
10
11 grid_search.fit(features_train,labels_train)
12
13
14 print(grid_search.grid_scores_)
15 print(grid_search.best_params_)
16 print(grid_search.best_score_)
17
```

/home/niklas/Documents/anaconda3/envs/python2/lib/python2.7/site-packages/sklearn/metrics/classification.py:1135: UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 due to no predicted samples.

'precision', 'predicted', average, warn_for)

/home/niklas/Documents/anaconda3/envs/python2/lib/python2.7/site-packages/sklearn/metrics/classification.py:1135: UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 due to no predicted samples.

'precision', 'predicted', average, warn_for)

/home/niklas/Documents/anaconda3/envs/python2/lib/python2.7/site-packages/sklearn/metrics/classification.py:1135: UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 due to no predicted samples.

'precision', 'predicted', average, warn_for)

/home/niklas/Documents/anaconda3/envs/python2/lib/python2.7/site-packages/sklearn/metrics/classification.py:1135: UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 due to no predicted samples.

'precision', 'predicted', average, warn_for)

/home/niklas/Documents/anaconda3/envs/python2/lib/python2.7/site-packages/sklearn/metrics/classification.py:1135: UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 due to no predicted samples.

'precision', 'predicted', average, warn_for)

/home/niklas/Documents/anaconda3/envs/python2/lib/python2.7/site-packages/sklearn/metrics/classification.py:1135: UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 due to no predicted samples.

'precision', 'predicted', average, warn_for)

/home/niklas/Documents/anaconda3/envs/python2/lib/python2.7/site-packages/sklearn/metrics/classification.py:1135: UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 due to no predicted samples.

'precision', 'predicted', average, warn_for)

/home/niklas/Documents/anaconda3/envs/python2/lib/python2.7/site-packages/sklearn/metrics/classification.py:1135: UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 due to no predicted samples.

'precision', 'predicted', average, warn_for)

/home/niklas/Documents/anaconda3/envs/python2/lib/python2.7/site-packages/sklearn/metrics/classification.py:1135: UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 due to no predicted samples.

'precision', 'predicted', average, warn_for)

/home/niklas/Documents/anaconda3/envs/python2/lib/python2.7/site-packages/sklearn/metrics/classification.py:1135: UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 due to no predicted samples.

'precision', 'predicted', average, warn_for)

/home/niklas/Documents/anaconda3/envs/python2/lib/python2.7/site-packages/sklearn/metrics/classification.py:1135: UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 due to no predicted samples.

'precision', 'predicted', average, warn_for)

/home/niklas/Documents/anaconda3/envs/python2/lib/python2.7/site-packages/sklearn/metrics/classification.py:1135: UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 due to no predicted samples.

'precision', 'predicted', average, warn_for)

[mean: 0.13069, std: 0.18856, params: {'classify__min_samples_split': 2, 'kbest__k': 6, 'classify__max_depth': 5}, mean: 0.15064, std: 0.11718, params: {'clas

```
In [195]: 1 score = grid_search.score(features_test, labels_test)
2 pred = grid_search.predict(features_test)
3 accuracy = accuracy_score(labels_test, pred)
4 prec = precision_score(labels_test, pred)
5 recall = recall_score(labels_test, pred)
6
7 print "RFC accuracy score:", "%.2f" % round(accuracy,3) , "precision:", "%.2f"
8
9 print confusion_matrix(labels_test, pred)
```

```
RFC accuracy score: 0.86 precision: 0.33 recall: 0.50
[[36  4]
 [ 2  2]]
```

usage of evaluation metrics

I have tuned the model in a way that it can correctly predict 2 of 4 POI of the test data set. (recall score) On the other hand, 2 persons were wrongly labeled as POI by the model, resulting in a precision score of 0.33.

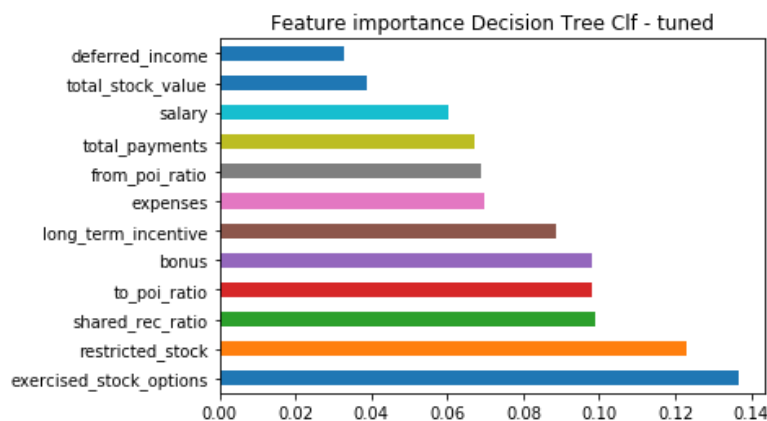
Discussion of parameter tuning

Hyperparameters are set before any Machine learning algorithm is run, hence, it becomes very essential to set an optimal value of hyperparameters as it effects the convergence of any algorithm to a large extent.

Comparing the feature importance plots before and after tuning indicates that this change had also an impact on the weight of the features. The most important feature still remains "exercised stock options" followed by "shared receipt ratio".

```
In [196]: 1 (pd.Series(clf.feature_importances_, index=features_train.columns)
2          .nlargest(12))
```

```
Out[196]: <matplotlib.axes._subplots.AxesSubplot at 0x7f5e45d1a650>
```



Running the tester.py script with this model gives following results:

```
Accuracy: 0.86553 Precision: 0.48791 Recall: 0.17150 F1: 0.25379 F2: 0.19706 Total predictions: 15000
True positives: 343 False positives: 360 False negatives: 1657 True negatives: 12640
```

This results is worse compared to the initial DecisionTreeClassifier with max_features=8

```
In [197]: 1 ### Task 6: Dump your classifier, dataset, and features_list so anyone can
2 ### check your results. You do not need to change anything below, but make s
3 ### that the version of poi_id.py that you submit can be run on its own and
4 ### generates the necessary .pkl files for validating your results.
5 clf = DecisionTreeClassifier(random_state=2, max_features=8)
6 clf.fit(features_train, labels_train)
7 dump_classifier_and_data(clf, my_dataset, features_list)
```

Final Results running the tester.py script

Accuracy: 0.81580 Precision: 0.31911 Recall: 0.33650 F1: 0.32757 F2: 0.33287 Total predictions: 15000

True positives: 673 False positives: 1436 False negatives: 1327 True negatives: 11564

In []:

1