**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"]**

Enron Corporation was U.S. energy-trading and utilities company founded in 1985. It achieved spectacular success in a very short period of time. This company claimed revenues of around 101 billion dollars in 2000. In only 15 years, this company was named America's Most innovative Company for six consecutive years by Fortune 500. Company stock price rose to 90 dollars. However, this success of Enron was short-lived. It started falling off more rapidly and bankruptcy was filed on Dec.2, 2001. Federal Energy Regulatory Commission investigated after the company's collapse and approximately 600,000 emails generated by employees of the Enron Corporation was made public which is called enron corpus.

In this project, main goal is to build a classification algorithm to predict a person of interest identifier (POI) based on email and financial features in the enron dataset involving 146 enron executives. A POI is anyone who is indicted, settled without admitting the guilt and testified in exchange for immunity. Predicted POI from this analysis and actual POI in the dataset will be compared to validate my prediction. Such model could be used to find additional suspects who were not indicted during the original investigation, or to find persons of interest during fraud investigations at other businesses. This study involves: -Exploration of Enron dataset -Selection of Features -Selection of Algorithms -Validation and Evaluation

```
In [1]:     1
            2  import matplotlib.pyplot as plt
            3  from IPython.display import Image
            4  import numpy as np
            5  import pandas as pd
            6  import sys
            7  import pickle
            8  from sklearn.naive_bayes import GaussianNB
            9  from sklearn.metrics import accuracy_score
           10  from sklearn.metrics import precision_score
           11  from sklearn.metrics import recall_score
           12  from sklearn.metrics import recall_score,precision_score
           13  from time import time
           14  from sklearn.grid_search import GridSearchCV
           15
           16  sys.path.append("../tools/")
           17
           18  from feature_format import featureFormat, targetFeatureSplit
           19
```

```
/anaconda/lib/python2.7/site-packages/sklearn/cross_validation.py:44:
DeprecationWarning: This module was deprecated in version 0.18 in favo
r of the model_selection module into which all the refactored classes
and functions are moved. Also note that the interface of the new CV it
erators are different from that of this module. This module will be re
moved in 0.20.
  "This module will be removed in 0.20.", DeprecationWarning)
/anaconda/lib/python2.7/site-packages/sklearn/grid_search.py:43: Depre
cationWarning: This module was deprecated in version 0.18 in favor of
the model_selection module into which all the refactored classes and f
unctions are moved. This module will be removed in 0.20.
  DeprecationWarning)
```

```
In [2]:     1  with open("final_project_dataset.pkl", "r") as data_file:
            2      data_dict = pickle.load(data_file)
```

```
In [3]:     1  enron_data = pd.DataFrame.from_dict(data_dict, orient = 'index')
```

## Exploration of Enron Dataset

Dataset was investigated. First structure of the dataset was examined.

In [4]:
```
1 enron_data.head()
```

Out[4]:

| | salary | to_messages | deferral_payments | total_payments | exercised_stock_options |
|---|---|---|---|---|---|
| **ALLEN PHILLIP K** | 201955 | 2902 | 2869717 | 4484442 | 1729541 |
| **BADUM JAMES P** | NaN | NaN | 178980 | 182466 | 257817 |
| **BANNANTINE JAMES M** | 477 | 566 | NaN | 916197 | 4046157 |
| **BAXTER JOHN C** | 267102 | NaN | 1295738 | 5634343 | 6680544 |
| **BAY FRANKLIN R** | 239671 | NaN | 260455 | 827696 | NaN |

5 rows × 21 columns

In [5]:
```
1 #structure of enron_data set
2 print "There were", len(enron_data), "persons involved."
```

There were 146 persons involved.

```
In [6]:    1  print data_dict.keys()
```

['METTS MARK', 'BAXTER JOHN C', 'ELLIOTT STEVEN', 'CORDES WILLIAM R',
'HANNON KEVIN P', 'MORDAUNT KRISTINA M', 'MEYER ROCKFORD G', 'MCMAHON
JEFFREY', 'HORTON STANLEY C', 'PIPER GREGORY F', 'HUMPHREY GENE E', 'U
MANOFF ADAM S', 'BLACHMAN JEREMY M', 'SUNDE MARTIN', 'GIBBS DANA R', '
LOWRY CHARLES P', 'COLWELL WESLEY', 'MULLER MARK S', 'JACKSON CHARLENE
R', 'WESTFAHL RICHARD K', 'WALTERS GARETH W', 'WALLS JR ROBERT H', 'KI
TCHEN LOUISE', 'CHAN RONNIE', 'BELFER ROBERT', 'SHANKMAN JEFFREY A', '
WODRASKA JOHN', 'BERGSIEKER RICHARD P', 'URQUHART JOHN A', 'BIBI PHILI
PPE A', 'RIEKER PAULA H', 'WHALEY DAVID A', 'BECK SALLY W', 'HAUG DAVI
D L', 'ECHOLS JOHN B', 'MENDELSOHN JOHN', 'HICKERSON GARY J', 'CLINE K
ENNETH W', 'LEWIS RICHARD', 'HAYES ROBERT E', 'MCCARTY DANNY J', 'KOPP
ER MICHAEL J', 'LEFF DANIEL P', 'LAVORATO JOHN J', 'BERBERIAN DAVID',
'DETMERING TIMOTHY J', 'WAKEHAM JOHN', 'POWERS WILLIAM', 'GOLD JOSEPH'
, 'BANNANTINE JAMES M', 'DUNCAN JOHN H', 'SHAPIRO RICHARD S', 'SHERRIF
F JOHN R', 'SHELBY REX', 'LEMAISTRE CHARLES', 'DEFFNER JOSEPH M', 'KIS
HKILL JOSEPH G', 'WHALLEY LAWRENCE G', 'MCCONNELL MICHAEL S', 'PIRO JI
M', 'DELAINEY DAVID W', 'SULLIVAN-SHAKLOVITZ COLLEEN', 'WROBEL BRUCE',
'LINDHOLM TOD A', 'MEYER JEROME J', 'LAY KENNETH L', 'BUTTS ROBERT H',
'OLSON CINDY K', 'MCDONALD REBECCA', 'CUMBERLAND MICHAEL S', 'GAHN ROB
ERT S', 'MCCLELLAN GEORGE', 'HERMANN ROBERT J', 'SCRIMSHAW MATTHEW', '
GATHMANN WILLIAM D', 'HAEDICKE MARK E', 'BOWEN JR RAYMOND M', 'GILLIS
JOHN', 'FITZGERALD JAY L', 'MORAN MICHAEL P', 'REDMOND BRIAN L', 'BAZE
LIDES PHILIP J', 'BELDEN TIMOTHY N', 'DURAN WILLIAM D', 'THORN TERENCE
H', 'FASTOW ANDREW S', 'FOY JOE', 'CALGER CHRISTOPHER F', 'RICE KENNET
H D', 'KAMINSKI WINCENTY J', 'LOCKHART EUGENE E', 'COX DAVID', 'OVERDY
KE JR JERE C', 'PEREIRA PAULO V. FERRAZ', 'STABLER FRANK', 'SKILLING J
EFFREY K', 'BLAKE JR. NORMAN P', 'SHERRICK JEFFREY B', 'PRENTICE JAMES
', 'GRAY RODNEY', 'PICKERING MARK R', 'THE TRAVEL AGENCY IN THE PARK',
'NOLES JAMES L', 'KEAN STEVEN J', 'TOTAL', 'FOWLER PEGGY', 'WASAFF GEO
RGE', 'WHITE JR THOMAS E', 'CHRISTODOULOU DIOMEDES', 'ALLEN PHILLIP K'
, 'SHARP VICTORIA T', 'JAEDICKE ROBERT', 'WINOKUR JR. HERBERT S', 'BRO
WN MICHAEL', 'BADUM JAMES P', 'HUGHES JAMES A', 'REYNOLDS LAWRENCE', '
DIMICHELE RICHARD G', 'BHATNAGAR SANJAY', 'CARTER REBECCA C', 'BUCHANA
N HAROLD G', 'YEAP SOON', 'MURRAY JULIA H', 'GARLAND C KEVIN', 'DODSON
KEITH', 'YEAGER F SCOTT', 'HIRKO JOSEPH', 'DIETRICH JANET R', 'DERRICK
JR. JAMES V', 'FREVERT MARK A', 'PAI LOU L', 'BAY FRANKLIN R', 'HAYSLE
TT RODERICK J', 'FUGH JOHN L', 'FALLON JAMES B', 'KOENIG MARK E', 'SAV
AGE FRANK', 'IZZO LAWRENCE L', 'TILNEY ELIZABETH A', 'MARTIN AMANDA K'
, 'BUY RICHARD B', 'GRAMM WENDY L', 'CAUSEY RICHARD A', 'TAYLOR MITCHE
LL S', 'DONAHUE JR JEFFREY M', 'GLISAN JR BEN F']

```
In [7]:    1 print data_dict["GLISAN JR BEN F"]
           2
           3
```

{'salary': 274975, 'to_messages': 873, 'deferral_payments': 'NaN', 'to
tal_payments': 1272284, 'exercised_stock_options': 384728, 'bonus': 60
0000, 'restricted_stock': 393818, 'shared_receipt_with_poi': 874, 'res
tricted_stock_deferred': 'NaN', 'total_stock_value': 778546, 'expenses
': 125978, 'loan_advances': 'NaN', 'from_messages': 16, 'other': 20030
8, 'from_this_person_to_poi': 6, 'poi': True, 'director_fees': 'NaN',
'deferred_income': 'NaN', 'long_term_incentive': 71023, 'email_address
': 'ben.glisan@enron.com', 'from_poi_to_this_person': 52}

```
In [8]:    1 (enron_data['poi'].value_counts()[True])
```

Out[8]: 18

```
1 enron_data.describe().transpose()
```

Out[9]:

| | count | unique | top | freq |
|---|---|---|---|---|
| salary | 146 | 95 | NaN | 51 |
| to_messages | 146 | 87 | NaN | 60 |
| deferral_payments | 146 | 40 | NaN | 107 |
| total_payments | 146 | 126 | NaN | 21 |
| exercised_stock_options | 146 | 102 | NaN | 44 |
| bonus | 146 | 42 | NaN | 64 |
| restricted_stock | 146 | 98 | NaN | 36 |
| shared_receipt_with_poi | 146 | 84 | NaN | 60 |
| restricted_stock_deferred | 146 | 19 | NaN | 128 |
| total_stock_value | 146 | 125 | NaN | 20 |
| expenses | 146 | 95 | NaN | 51 |
| loan_advances | 146 | 5 | NaN | 142 |
| from_messages | 146 | 65 | NaN | 60 |
| other | 146 | 93 | NaN | 53 |
| from_this_person_to_poi | 146 | 42 | NaN | 60 |
| poi | 146 | 2 | False | 128 |
| director_fees | 146 | 18 | NaN | 129 |
| deferred_income | 146 | 45 | NaN | 97 |
| long_term_incentive | 146 | 53 | NaN | 80 |
| email_address | 146 | 112 | NaN | 35 |
| from_poi_to_this_person | 146 | 58 | NaN | 60 |

Enron dataset consists of lot of missing values (NaN). All the features have at least one missing values. Some features have more than 50% of their values missing, as shown by the table above. NaNs are replaced by 0 for training our algorithm later.
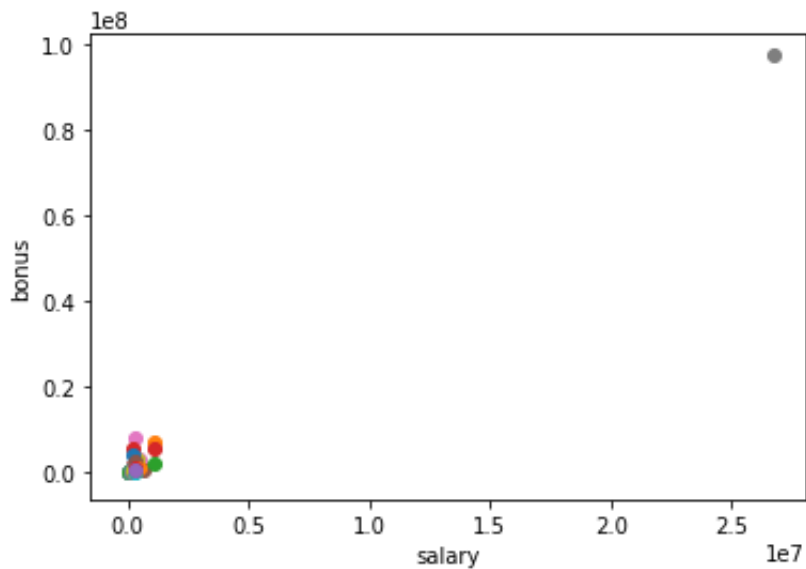
In [10]:

```
1 enron_data.replace(to_replace='NaN', value=0.0, inplace=True)
```

Scatter plots are great tools for finding outliers. So first I plotted salary against bonus in a scatter plot and looked for the data distribution.

```
In [11]:    1  from IPython.display import Image
            2  features = ["salary", "bonus"]
            3  #data_dict.pop('TOTAL', 0)
            4  data = featureFormat(data_dict, features)
            5  ### plot features
            6  for point in data:
            7      salary = point[0]
            8      bonus = point[1]
            9      plt.scatter( salary, bonus )
           10
           11  plt.xlabel("salary")
           12  plt.ylabel("bonus")
           13  plt.show()
```



The above plot allowed to see an outlier named 'Total'which was sum of all the data points in the given plot(artifact in scatterplot). So this outlier was excluded from our analysis manually. There were several other high number values (outliers) which could possibly represent POIs, thus were included in the dataset. Scatter plot was drawn after removing outliers as below.

```
In [12]:   1  ##2 removing outliers
           2  features = ["salary", "bonus"]
           3  data_dict.pop('TOTAL', 0)
           4  data = featureFormat(data_dict, features)
           5
           6  #remove NAN from dataset
           7  outliers = []
           8  for key in data_dict:
           9      val = data_dict[key]['salary']
          10      if val == 'NaN':
          11          continue
          12      outliers.append((key, int(val)))
          13  outliers_final = (sorted(outliers, key=lambda x:x[1], reverse = True)
          14  print outliers_final
```
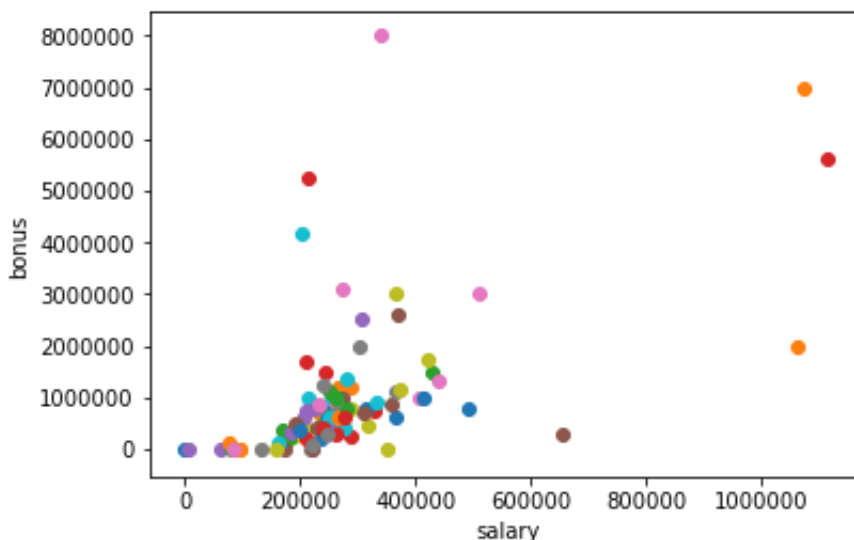
[('SKILLING JEFFREY K', 1111258), ('LAY KENNETH L', 1072321), ('FREVER
T MARK A', 1060932), ('PICKERING MARK R', 655037), ('WHALLEY LAWRENCE
G', 510364), ('DERRICK JR. JAMES V', 492375), ('FASTOW ANDREW S', 4406
98), ('SHERRIFF JOHN R', 428780), ('RICE KENNETH D', 420636), ('CAUSEY
RICHARD A', 415189)]

```
In [13]:   1  #scatterplot after outlier removal
           2
           3  features = ["salary", "bonus"]
           4  #data_dict.pop('TOTAL', 0)
           5  data = featureFormat(data_dict, features)
           6  ### plot features
           7  for point in data:
           8      salary = point[0]
           9      bonus = point[1]
          10      plt.scatter( salary, bonus )
          11
          12  plt.xlabel("salary")
          13  plt.ylabel("bonus")
          14  plt.show()
```

# Selection of Feature/Feature engineering

**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 importances 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", "intelligently select features", "properly scale features"]**

Feature selection in machine learning is selection of a feature which potentially would show some sort of pattern in the prediction analysis. Enron dataset consists two set of features - financial(salary, bonus, stock etc) and communication(to and from emails). Here I am more interested in emails those involve POIs. The assumption is that "Communication between POI and POI would be more frequent than between POI and non-POIs". I attempted to create two features, fractions of emails this person gets from POI(fraction_from_poi) and fraction of emails this person sends to POI(fraction_to_poi).

fraction_from_poi = number of emails this person gets from POI/total number of from messages

fraction_to_poi = number of emails this person sends to POI/total number sent messages

If email data is 'NaN', ratio is set to 0.

In [14]:
```python
#3 create new features
#new features: fraction_to_poi = fraction of emails sent to POIs, fra
def dict_to_list(key,normalizer):
    new_list=[]

    for i in data_dict:
        if data_dict[i][key]=="NaN" or data_dict[i][normalizer]=="NaN
            new_list.append(0.)
        elif data_dict[i][key]>=0:
            new_list.append(float(data_dict[i][key])/float(data_dict[
    return new_list

### create two lists of new features
fraction_from_poi = dict_to_list("from_poi_to_this_person","to_messag
fraction_to_poi = dict_to_list("from_this_person_to_poi","from_messag

### insert new features into data_dict
count=0
for i in data_dict:
    data_dict[i]["fraction_from_poi"]=fraction_from_poi[count]
```

```
21        data_dict[i]["fraction_to_poi"]=fraction_to_poi[count]
22        count +=1
23
24
25 features_list = ["poi", "fraction_from_poi", "fraction_to_poi"]
26        ### store to my_dataset for easy export below
27 my_dataset = data_dict
28
29
30 ### these two lines extract the features specified in features_list
31 ### and extract them from data_dict, returning a numpy array
32 data = featureFormat(my_dataset, features_list)
33
34 ### plot new features
35 for point in data:
36        from_poi = point[1]
37        to_poi = point[2]
38        plt.scatter( from_poi, to_poi )
39        if point[0] == 1:
40            plt.scatter(from_poi, to_poi, color="r", marker="*")
41 plt.xlabel("fraction of emails this person gets from poi")
42 plt.show()
43
```



In [15]:
```
1 print features_list
```

['poi', 'fraction_from_poi', 'fraction_to_poi']

To select more impactful features for classification, 'feature_importances' attribute of "Decision Tree" was used. Features were ranked and selection process was half manual process. First all the possible features were included in features_list and then selection was done on the basis of feature ranking.

In [16]:
```
1
2
3  features list — ["poi"    "salary"    "bonus"    "fraction from poi"    "frac
```

```
 3  features_list = [ poi , salary , bonus , fraction_from_poi , frac
 4                     'deferral_payments', 'total_payments', 'loan_advance
 5                     'deferred_income', 'total_stock_value', 'expenses',
 6                     'long_term_incentive', 'shared_receipt_with_poi', 'r
 7  data = featureFormat(my_dataset, features_list)
 8
 9  ### split into labels and features (this line assumes that the first
10  ### feature in the array is the label, which is why "poi" must always
11  ### be first in features_list
12  labels, features = targetFeatureSplit(data)
13
14  ### split data into training and testing datasets
15  #deploying feature selection
16
17  from sklearn import cross_validation
18  from sklearn.cross_validation import train_test_split
19  features_train, features_test, labels_train, labels_test = cross_vali
20
21
22
23  ##try Decision tree
24  from sklearn.tree import DecisionTreeClassifier
25
26  t0 = time()
27
28  clf = DecisionTreeClassifier()
29  clf.fit(features_train,labels_train)
30  score = clf.score(features_test,labels_test)
31  print 'Decision Tree'
32  print 'accuracy before tuning', score
33  print "Decision tree algorithm time:", round(time()-t0, 3), "s"
34
35  importances = clf.feature_importances_
36  import numpy as np
37  indices = np.argsort(importances)[::-1]
38  print 'Feature Ranking: '
39  for i in range(16):
40      print "{} feature {} ({})".format(i+1,features_list[i+1],importan
41
42
43
44
45
46
47
48
```

```
Decision Tree
accuracy before tuning 0.666666666667
Decision tree algorithm time: 0.007 s
Feature Ranking:
1 feature salary (0.174317879326)
2 feature bonus (0.158290984378)
3 feature fraction_from_poi (0.14622972935)
```

```
4 feature fraction_to_poi (0.128198757764)
5 feature deferral_payments (0.118337314859)
6 feature total_payments (0.0955181169023)
7 feature loan_advances (0.0879795396419)
8 feature restricted_stock_deferred (0.0534161490683)
9 feature deferred_income (0.0377115287109)
10 feature total_stock_value (0.0)
11 feature expenses (0.0)
12 feature exercised_stock_options (0.0)
13 feature long_term_incentive (0.0)
14 feature shared_receipt_with_poi (0.0)
15 feature restricted_stock (0.0)
16 feature director_fees (0.0)
```

We select and keep least number of features which hold maximum information and show pattern and trends in data. From here I selected 9 features: ["salary", "bonus", "fraction_from_poi", "fraction_to_poi", "deferral_payments", "total_payments", "loan_advances", "restricted_stock_deferred", "deferred_income"] Accuracy for this feature set was 0.73. Recall and precision scores were too low, so I manually chose features which gave recall and precision score higher than 0.3. My final feature selections were: ["fraction_from_poi", "fraction_to_poi","shared_receipt_with_poi"]

# Algorithm selection and tuning

```
In [17]:  1
          2  features_list = ["poi", "fraction_from_poi", "fraction_to_poi","share
          3
          4
          5  ##try Naive Bayes
          6
          7  t0 = time()
          8  clf = GaussianNB()
          9  clf.fit(features_train, labels_train)
         10  pred = clf.predict(features_test)
         11  print "Naive Bayes recall score", (recall_score(labels_test,pred))
         12  print "Naive Bayes precision score", (precision_score(labels_test,pre
         13  #print accuracy #(clf.score(features_test, labels_test))
         14  accuracy = accuracy_score(pred,labels_test)
         15  print 'accuracy', accuracy
         16  print "NB algorithm time:", round(time()-t0, 3), 's'
         17
         18
         19
         20  ###Adaboost
         21  from sklearn.ensemble import AdaBoostClassifier
         22  clf = AdaBoostClassifier(DecisionTreeClassifier(min_samples_split = 4
         23                           algorithm="SAMME",
         24                           n_estimators=200)
         25  clf.fit(features_train,labels_train)
         26  pred = clf.predict(features_test)
         27
         28  acc =  accuracy_score(pred,labels_test)
         29  print "ADABOOST:"
         30  print acc
         31  print "AB Recall Score" + str(recall_score(labels_test, pred))
         32  print "AB Precision Score" + str(precision_score(labels_test, pred))
         33
         34
         35
```

```
Naive Bayes recall score 0.5
Naive Bayes precision score 0.181818181818
accuracy 0.266666666667
NB algorithm time: 0.01 s
ADABOOST:
0.8
AB Recall Score0.25
AB Precision Score1.0
```

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

**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? What parameters did you tune?**

Naive Bayes (accuracy = 0.27), decision tree (accuracy = 0.67) and adaboost(accuracy= 0.73) algorithms were applied. Accuracy was lower with Naive Bayes, could probably due to distribution pattern of features in the dataset. Adaboost showed decent accuracy but had recall score of less than 0.3. Here, Decision tree would be more appropriate algorithm for POI identifier, since as it has accuracy of 0.67 before any tuning. Besides it is more efficient in finding irregular decision boundaries and does not need feature scaling. To optimize its performance, parameters like min_sample_split could be varied.

According to the enron dataset there were only 18 POIs. Since there were few POI subjects in the dataset, precision and recall score were considered better evaluater. I chose decision tree as a final algorithm. After final algorithm was chosen, algorithm parameter min_samples_split was tuned manually. We varied the min_samples_split from 2 to 7 and compared the recall and precision score. We found recall and precision score did not change between numbers 3 to 7.

In [18]:

```
1  ### use manual tuning parameter min_samples_split
2  t0 = time()
3  clf = DecisionTreeClassifier(min_samples_split=5)
4  clf = clf.fit(features_train,labels_train)
5  pred= clf.predict(features_test)
6  print("done in %0.3fs" % (time() - t0))
7
8  acc=accuracy_score(labels_test, pred)
9
10 print "Validating algorithm:"
11 print "accuracy after tuning = ", acc
12
13 # function for calculation ratio of true positives
14 # out of all positives (true + false)
15 print 'precision = ', precision_score(labels_test,pred)
16
17 # function for calculation ratio of true positives
18 # out of true positives and false negatives
19 print 'recall = ', recall_score(labels_test,pred)
20
21 ### dump your classifier, dataset and features_list so
22 ### anyone can run/check your results
23 pickle.dump(clf, open("my_classifier.pkl", "w") )
24 pickle.dump(data_dict, open("my_dataset.pkl", "w") )
25 pickle.dump(features_list, open("my_feature_list.pkl", "w") )
26
```

```
done in 0.006s
Validating algorithm:
accuracy after tuning =  0.733333333333
precision =  0.5
recall =  0.5
```

```
In [19]:   1  #t0 = time()
           2  #clf = DecisionTreeClassifier(min_samples_split=2)
           3  #clf = clf.fit(features_train,labels_train)
           4  #pred= clf.predict(features_test)
           5  #print("done in %0.3fs" % (time() - t0))
           6
           7  #acc=accuracy_score(labels_test, pred)
           8
           9  #print "Validating algorithm:"
          10  #print "accuracy after tuning = ", acc
          11
          12  #function for calculation ratio of true positives
          13  #out of all positives (true + false)
          14  #print 'precision = ', precision_score(labels_test,pred)
          15
          16  #function for calculation ratio of true positives
          17  #out of true positives and false negatives
          18  #print 'recall = ', recall_score(labels_test,pred)
```

```
 1  min_samples_split      precision      recall
 2          2                  0.4          0.67
 3          3                  0.5          0.67
 4          4                  0.5          0.67
 5          5                  0.67         0.67
 6          6                  0.67         0.67
 7          7                  0.67         0.67
 8
```

## Analysis Validation and Performance

### 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 the process where we determine the robustness of our predictive models. In our case, analysis was validated using K-fold validation. Such validation process enhances the likelihood that our algorithm will be reliable and robust. A classic mistake in validation process is called over-fitting, where the model is trained and it performs very well on the training dataset, but is actually worse on the cross-validation and test datasets.

```
In [20]:   1  ### features_list is a list of strings, each of which is a feature na
           2  ### first feature must be "poi", as this will be singled out as the l
           3  features_list = ["poi", "fraction_from_poi", "fraction_to_poi", "shar
           4
           5
           6  ### store to my_dataset for easy export below
           7  my_dataset = data_dict
           8
           9
          10  ### these two lines extract the features specified in features_list
```

```python
11 ### and extract them from data_dict, returning a numpy array
12 data = featureFormat(my_dataset, features_list)
13
14
15 ### split into labels and features (this line assumes that the first
16 ### feature in the array is the label, which is why "poi" must always
17 ### be first in features_list
18 labels, features = targetFeatureSplit(data)
19
20
21 ### machine learning goes here!
22 ### please name your classifier clf for easy export below
23
24 ### deploying feature selection
25 from sklearn import cross_validation
26 features_train, features_test, labels_train, labels_test = cross_vali
27
28 ### use KFold for split and validate algorithm
29 from sklearn.cross_validation import KFold
30 kf=KFold(len(labels),3)
31 for train_indices, test_indices in kf:
32     #make training and testing sets
33     features_train= [features[ii] for ii in train_indices]
34     features_test= [features[ii] for ii in test_indices]
35     labels_train=[labels[ii] for ii in train_indices]
36     labels_test=[labels[ii] for ii in test_indices]
37
38 from sklearn.tree import DecisionTreeClassifier
39
40 t0 = time()
41
42 clf = DecisionTreeClassifier()
43 clf.fit(features_train,labels_train)
44 score = clf.score(features_test,labels_test)
45 print 'accuracy before tuning ', score
46
47 print "Decision tree algorithm time:", round(time()-t0, 3), "s"
48
49
50 ### use manual tuning parameter min_samples_split
51 t0 = time()
52 clf = DecisionTreeClassifier(min_samples_split=5)
53 clf = clf.fit(features_train,labels_train)
54 pred= clf.predict(features_test)
55 print("done in %0.3fs" % (time() - t0))
56
57 acc=accuracy_score(labels_test, pred)
58
59 print "Validating algorithm:"
60 print "accuracy after tuning = ", acc
61
62 # function for calculation ratio of true positives
63 # out of all positives (true + false)
64 print 'precision = ', precision_score(labels_test,pred)
```

```
65
66  # function for calculation ratio of true positives
67  # out of true positives and false negatives
68  print 'recall = ', recall_score(labels_test,pred)
69
70
71  ### dump your classifier, dataset and features_list so
72  ### anyone can run/check your results
73  pickle.dump(clf, open("my_classifier.pkl", "w") )
74  pickle.dump(data_dict, open("my_dataset.pkl", "w") )
75  pickle.dump(features_list, open("my_feature_list.pkl", "w") )
```

```
accuracy before tuning  0.857142857143
Decision tree algorithm time: 0.003 s
done in 0.001s
Validating algorithm:
accuracy after tuning =  0.892857142857
precision =  0.5
recall =  0.666666666667
```

# Discussion and Conclusions

As this dataset was small and imbalanced dataset, accuracy was not a good metric for evaluating the algorithm. So, we used different scoring metric; precision and recall score. Decision Tree classifier showed the precision score and recall score higher than 0.3. Comparatively, Adaboost and Naive Bayes showed low accuracy and low precision score.

**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 chose Decision Tree as a final algorithm to predict whether POI identified through this test is indeed POI. Feature scaling was not done as it is not required while using Decision Tree. The precision score is the probability that the person identified as POI is infact POI. In our case, precision = 0.67 means 67 percent of the time POI identified in the test would be the real POI whereas 33% of the time this test could flag a wrong person. These numbers could be increased by changing or exploring more into email information. Finally, Our prediction was validated by using K-fold validation.

Thus in this project, I used machine learning algorithms to identify POI in the Enron dataset. Naive Bayes, Adaboost and Decision Tree algorithms were compared. Decision Tree showed best precision and recall score among them showing 0.67 and 0.67 respectively. Small size of data set and even smaller number of POI made this analysis more challanging.

# References:

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

https://www.cs.cmu.edu/~enron/ (https://www.cs.cmu.edu/~enron/)

https://www.technologyreview.com/s/515801/the-immortal-life-of-the-enron-e-mails/
(https://www.technologyreview.com/s/515801/the-immortal-life-of-the-enron-e-mails/)

https://classroom.udacity.com/nanodegrees/nd002/parts/0021345409/modules/3174288624754617
(https://classroom.udacity.com/nanodegrees/nd002/parts/0021345409/modules/31742886247546

https://discussions.udacity.com/c/nd002-intro-to-machine-learning
(https://discussions.udacity.com/c/nd002-intro-to-machine-learning)

In [ ]: 1