

## PROJECT

# Identify Fraud from Enron Email

A part of the Data Analyst Nanodegree Program

## PROJECT REVIEW

## CODE REVIEW 4

## NOTES

▼ poi\_id.py 4

```
1 #!/usr/bin/python
```

AWESOME

Code is well structured with functions well commented! 🍌

```
2
3 '''
4 Harold Finch: [Opening narration from Season One]
5 You are being watched. The government has a secret system, a machine that spies on you every hour of every day.
6 I know because I built it. I designed the machine to detect acts of terror but it sees everything.
7 Violent crimes involving ordinary people, people like you. Crimes the government considered "irrelevant."
8 They wouldn't act, so I decided I would. But I needed a partner, someone with the skills to intervene.
9 Hunted by the authorities, we work in secret. You'll never find us, but victim or perpetrator, if your number's up... we'll find *you*.
```

AWESOME



```
10 '''
11
12 import matplotlib.pyplot as plt
13 import numpy as np
14 import pandas
15 import seaborn as sns
16 import math
17
18 from sklearn.preprocessing import MinMaxScaler
19 from sklearn.cross_validation import train_test_split
20 from sklearn.cross_validation import StratifiedShuffleSplit
21 from sklearn.cross_validation import StratifiedKFold
22 from sklearn.decomposition import PCA
23
24 from sklearn.svm import SVC
25 from sklearn.naive_bayes import GaussianNB
26 from sklearn.tree import DecisionTreeClassifier
27 from sklearn.cluster import KMeans
28 from sklearn.ensemble import AdaBoostClassifier
29
30 from sklearn.pipeline import Pipeline
31 from sklearn.model_selection import GridSearchCV
32 from sklearn.model_selection import RandomizedSearchCV
33 from sklearn.metrics import classification_report
34
35 plt.interactive(False)
36
37
38 sequence_nb = 0
39
```

```

40 #### Global variables to adjust the software behaviour
41 showunivariate = True
42 showheatmap = True
43 showcorrelation = True
44 performalltunings = True
45
46 '''
47     -----
48     Global tools
49     -----
50 '''
51
52 def GetSequence():
53     # This function return a sequence number
54     # It will be used to uniquely identify a season chart
55     global sequence_nb
56     sequence_nb = sequence_nb + 1
57     return sequence_nb
58
59 def GetValues(data, feature1, feature2=""):
60     # This function return one or two features from the project data dictionary
61     # Only lines with values are extracted
62     # When two features are mentioned, only the lines with the two valued
63     # (not NaN) features are provided
64
65     values1 = []
66     values2 = []
67     for item in data:
68         if feature2 == "":
69             if data[item][feature1] != 'NaN':
70                 values1.append(data[item][feature1])
71         else:
72             if data[item][feature1] != 'NaN' and data[item][feature2] != 'NaN':
73                 values1.append(data[item][feature1])
74                 values2.append(data[item][feature2])
75
76     if feature2 == "":
77         return values1
78     else:
79         return values1, values2
80
81 def ShowHist(data, features):
82     # This function display an histogram for all provided features
83
84     for feature in features:
85         values = GetValues(data, feature)
86         sns.plt.figure(GetSequence())
87         sns.distplot(values, axlabel = feature)
88         sns.plt.show()
89
90 def ShowScatter(data, feature1, feature2):
91     # This function display an scatter plot for two given features
92
93     values1, values2 = GetValues(data, feature1, feature2)
94
95     df = pandas.concat([pandas.Series(values1, name=feature1), pandas.Series(values2, name=feature2)], axis=1)
96     sns.plt.figure(GetSequence())
97     sns.regplot(x=feature1, y=feature2, data=df, label = feature1 + " vs " + feature2)
98     sns.plt.show()
99
100
101
102 def removeOutliers(data, feature, removeLower = True, removeUpper = True):
103     # Remove outliers for one specific feature
104     # It uses the interquantile method to identify outliers
105
106     values = []
107     nbRemoved = 0
108
109     # Extract all values for this features
110     values = GetValues(data, feature)
111     df = pandas.DataFrame(values)
112
113     Q1 = df.quantile(0.25)
114     Q3 = df.quantile(0.75)
115     IQR = Q3 - Q1
116
117     # Now, remove outliers
118     for item in data:
119         if feature in data[item]:
120             if data[item][feature] != 'NaN':
121                 if removeLower and data[item][feature] < Q1 - 1.5*IQR:
122                     del data[item][feature]
123                     data[item][feature] = 'NaN'
124                     nbRemoved = nbRemoved + 1
125                 if data[item][feature] != 'NaN':
126                     if removeUpper and data[item][feature] > Q3 + 1.5*IQR:
127                         del data[item][feature]
128                         data[item][feature] = 'NaN'
129                         nbRemoved = nbRemoved + 1
130     print "Outliers, feature ", feature, " - ", nbRemoved, " values removed"
131
132
133 def removeAllOutliers(data, featureList):
134     # This function removes outliers of all provided features
135
136     for feature in featureList:
137         removeOutliers(data, feature)

```

```

138
139 def RemoveFeature(data, feature):
140     # Remove one specific feature from the data
141
142     for item in data:
143         if feature in data[item]:
144             del data[item][feature]
145
146 def RemovePositiveNegative(data, feature, removeNegative=False, removePositive=False):
147     # Remove one specific feature value if positive and/or negative
148     # by default, it does nothing
149
150     nbRemoved = 0
151     for item in data:
152         if feature in data[item]:
153             if removeNegative and data[item][feature] != 'NaN':
154                 if data[item][feature] < 0:
155                     del data[item][feature]
156                     data[item][feature] = 'NaN'
157                     nbRemoved = nbRemoved + 1
158             if removePositive and data[item][feature] != 'NaN':
159                 if data[item][feature] > 0:
160                     del data[item][feature]
161                     data[item][feature] = 'NaN'
162                     nbRemoved = nbRemoved + 1
163
164     print "RemovePositiveNegative for feature", feature, " : ", nbRemoved, " removed values"
165
166 def TransformLog(data, feature, positive=True):
167
168     # This function transform a given feature into its log value
169     for item in data:
170         if feature in data[item]:
171             if data[item][feature] != 'NaN':
172                 value = data[item][feature]
173                 if positive and value >= 0:
174                     del data[item][feature]
175                     data[item][feature] = math.log10(1+value)
176                 else:
177                     del data[item][feature]
178                     data[item][feature] = math.log10(1+abs(value))
179
180 def MakeLog(data_dict):
181     TransformLog(data_dict, 'bonus', positive=True)
182     TransformLog(data_dict, 'long_term_incentive', positive=True)
183     TransformLog(data_dict, 'deferred_income', positive=False)
184     TransformLog(data_dict, 'deferral_payments', positive=True)
185     TransformLog(data_dict, 'other', positive=True)
186     TransformLog(data_dict, 'expenses', positive=True)
187
188     TransformLog(data_dict, 'total_payments', positive=True)
189     TransformLog(data_dict, 'exercised_stock_options', positive=True)
190     TransformLog(data_dict, 'restricted_stock', positive=True)
191     TransformLog(data_dict, 'total_stock_value', positive=True)
192
193     TransformLog(data_dict, 'to_messages', positive=True)
194     TransformLog(data_dict, 'from_poi_to_this_person', positive=True)
195     TransformLog(data_dict, 'from_messages', positive=True)
196     TransformLog(data_dict, 'from_this_person_to_poi', positive=True)
197     TransformLog(data_dict, 'shared_receipt_with_poi', positive=True)
198
199 def CreateRatio(data):
200     # Create the two new ratios
201
202     for item in data:
203         to_messages = data[item]['to_messages']
204         from_poi_to_this_person = data[item]['from_poi_to_this_person']
205         from_messages = data[item]['from_messages']
206         from_this_person_to_poi = data[item]['from_this_person_to_poi']
207
208         if to_messages == 'NaN' or from_poi_to_this_person == 'NaN':
209             data[item]['poi_to_ratio'] = 'NaN'
210         else:
211             data[item]['poi_to_ratio'] = float(from_poi_to_this_person) / float(to_messages)
212
213         if from_messages == 'NaN' or from_this_person_to_poi == 'NaN':
214             data[item]['poi_from_ratio'] = 'NaN'
215         else:
216             data[item]['poi_from_ratio'] = float(from_this_person_to_poi) / float(from_messages)
217
218     print "New Ratio crated"
219
220 '''
221
222     -----
223     Machine learning algorithm tuning
224     -----
225 '''
226
227 def test_classifier_optim(clf, dataset, feature_list, folds = 1000):
228     labels, features = targetFeatureSplit(dataset)
229     cv = StratifiedShuffleSplit(labels, folds, random_state = 42)
230     true_negatives = 0
231     false_negatives = 0
232     true_positives = 0
233     false_positives = 0
234     for train_idx, test_idx in cv:
235         features_train = []

```

```

236     features_test = []
237     labels_train = []
238     labels_test = []
239     for ii in train_idx:
240         features_train.append( features[ii] )
241         labels_train.append( labels[ii] )
242     for jj in test_idx:
243         features_test.append( features[jj] )
244         labels_test.append( labels[jj] )
245
246     ### fit the classifier using training set, and test on test set
247     clf.fit(features_train, labels_train)
248     predictions = clf.predict(features_test)
249     for prediction, truth in zip(predictions, labels_test):
250         if prediction == 0 and truth == 0:
251             true_negatives += 1
252         elif prediction == 0 and truth == 1:
253             false_negatives += 1
254         elif prediction == 1 and truth == 0:
255             false_positives += 1
256         elif prediction == 1 and truth == 1:
257             true_positives += 1
258         else:
259             print "Warning: Found a predicted label not == 0 or 1."
260             print "All predictions should take value 0 or 1."
261             print "Evaluating performance for processed predictions:"
262             break
263     try:
264         total_predictions = true_negatives + false_negatives + false_positives + true_positives
265         accuracy = 1.0*(true_positives + true_negatives)/total_predictions
266         precision = 1.0*true_positives/(true_positives+false_positives)
267         recall = 1.0*true_positives/(true_positives+false_negatives)
268         f1 = 2.0 * true_positives/(2*true_positives + false_positives+false_negatives)
269         f2 = (1+2.0*2.0) * precision*recall/(4*precision + recall)
270         return f1
271     except:
272         print "Got a divide by zero when trying out:", clf
273         print "Precision or recall may be undefined due to a lack of true positive predicitions."
274         return 0
275
276
277 def TuneGNB(data, features_list, max_features, verbose=0, val_strategy='basic'):
278
279     print "Tuning Gaussian Naive Bayes"
280
281     labels, features = targetFeatureSplit(data)
282     features_train, features_test, labels_train, labels_test = train_test_split(features, labels, test_size=0.3, random_state=42)
283     print "Training test size = ", len(features_train)
284
285     scaler = MinMaxScaler(feature_range=(0, 1))
286     pca = PCA(iterated_power='auto', n_components=None, random_state=None, svd_solver='auto', tol=0.0, whiten=False)
287     classifier = GaussianNB()
288
289     pipe = Pipeline(steps=[('scaler', scaler), ('pca', pca), ('gnb', classifier)])

```

AWESOME

Nice work building your pipeline!

```

290
291     params = dict(
292         pca_n_components=range(3, max_features))
293     if val_strategy == 'basic':
294         estimator = GridSearchCV(pipe, params, scoring='f1', verbose=verbose)

```

AWESOME

Excellent optimization of GridCV!

```

295     else:
296         # kfold
297         cv = StratifiedKFold(labels_train, 10)
298         estimator = GridSearchCV(pipe, params, scoring='f1', verbose=verbose, cv=cv)
299
300     estimator.fit(features_train, labels_train)
301
302     # Return the best number of components and associated score
303     clf = estimator.best_estimator_
304     f1 = test_classifier_optim(clf, data, features_list)
305
306     print "Best f1 vs best estimator score = ", f1, estimator.best_score_
307
308     print "Estimator parameters"
309     print "PCA number of components = ", clf.named_steps['pca'].n_components
310
311     return clf, f1
312
313
314 def TuneSVM(data, features_list, max_features, verbose=0, val_strategy='basic'):
315
316     # This function will try to find the best principal component decomposition
317     # I use a Support Vector Machine to assess decomposition performance
318

```

```

319 print "Tuning Support Vector Machine"
320
321 labels, features = targetFeatureSplit(data)
322 features_train, features_test, labels_train, labels_test = train_test_split(features, labels, test_size=0.3, random_state=42)
323 print "Training test size = ", len(features_train)
324
325 scaler = MinMaxScaler(feature_range=(0, 1))
326 pca = PCA(iterated_power='auto', n_components=None, random_state=None, svd_solver='auto', tol=0.0, whiten=False)
327 classifier = SVC(kernel='rbf', C=10000)
328
329 pipe = Pipeline(steps=[('scaler', scaler), ('pca', pca), ('svm', classifier)])
330
331 params = dict(
332     pca__n_components=range(3, max_features),
333     svm__C=[1, 10, 100, 1000],
334     svm__gamma=[0.01, 0.001, 0.0001],
335     svm__kernel=['rbf', 'linear', 'poly'])
336 if val_strategy == 'basic':
337     estimator = RandomizedSearchCV(pipe, params, scoring='f1', verbose=verbose, random_state=42)
338 else:
339     # kfold
340     cv = StratifiedKFold(labels_train, 10)
341     estimator = RandomizedSearchCV(pipe, params, scoring='f1', verbose=verbose, cv=cv, random_state=42)
342
343 estimator.fit(features_train, labels_train)
344
345 # Return the best number of components and associated score
346 clf = estimator.best_estimator_
347 f1 = test_classifier_optim(clf, data, features_list)
348 print "Best f1 vs best estimator score = ", f1, estimator.best_score_
349
350 print "Estimator parameters"
351 print "PCA number of components = ", clf.named_steps['pca'].n_components
352 print "SVM C = ", clf.named_steps['svm'].C
353 print "SVM gamma = ", clf.named_steps['svm'].gamma
354 print "SVM kernel = ", clf.named_steps['svm'].kernel
355
356 return clf, f1
357
358 def TuneDT(data, features_list, max_features, verbose=0, val_strategy='basic'):
359
360     print "Tuning Decision Tree"
361
362     labels, features = targetFeatureSplit(data)
363     features_train, features_test, labels_train, labels_test = train_test_split(features, labels, test_size=0.3, random_state=42)
364     print "Training test size = ", len(features_train)
365
366     scaler = MinMaxScaler(feature_range=(0, 1))
367     pca = PCA(iterated_power='auto', n_components=None, random_state=None, svd_solver='auto', tol=0.0, whiten=False)
368     classifier = DecisionTreeClassifier()
369
370     pipe = Pipeline(steps=[('scaler', scaler), ('pca', pca), ('dt', classifier)])
371
372     params = dict(
373         pca__n_components=range(3, max_features),
374         dt__criterion=['gini', 'entropy'],
375         dt__max_features=['sqrt', 'log2', None])
376
377     if val_strategy == 'basic':
378         estimator = RandomizedSearchCV(pipe, params, scoring='f1', verbose=verbose, random_state=42)
379     else:
380         # kfold
381         cv = StratifiedKFold(labels_train, 10)
382         estimator = RandomizedSearchCV(pipe, params, scoring='f1', verbose=verbose, cv=cv, random_state=42)
383
384     estimator.fit(features_train, labels_train)
385
386     # Return the best number of components
387     clf = estimator.best_estimator_
388     f1 = test_classifier_optim(clf, data, features_list)
389     print "Best f1 vs best estimator score = ", f1, estimator.best_score_
390
391     print "Estimator parameters"
392     print "PCA number of components = ", clf.named_steps['pca'].n_components
393     print "DT criterion = ", clf.named_steps['dt'].criterion
394     print "DT max features = ", clf.named_steps['dt'].max_features
395
396     return clf, f1
397
398
399
400
401 import sys
402 import pickle
403 sys.path.append("../tools/")
404
405 from feature_format import featureFormat, targetFeatureSplit
406 from tester import dump_classifier_and_data
407
408
409
410 ### Load the dictionary containing the dataset
411 with open("final_project_dataset.pkl", "r") as data_file:
412     data_dict = pickle.load(data_file)
413
414 # Task 1: Get a global overview on the data
415
416 print "Descriptive statistics of salary"

```

```

418 df = pandas.DataFrame(featureFormat(data_dict, ['salary'], sort_keys = False,remove_NaN=False),columns=['salary'])
419 print df.describe()
420 # Delete TOTAL value and perform summary statistics on all features
421
422 print "TOTAL line deletion"
423 del data_dict['TOTAL']
424
425 print "Descriptive statistics on all features"
426 allFeatures = ['salary', 'bonus', 'long_term_incentive', 'deferred_income', 'deferral_payments', 'loan_advances', 'other', 'expenses', 'dir
427 df = pandas.DataFrame(featureFormat(data_dict, allFeatures, sort_keys = False,remove_NaN=False),columns=allFeatures)
428 print df.describe()
429
430 # Just let see how many poi and non poi we have.
431 pois = df['poi']
432 print "POIs = ",pois.sum(), " / Non POIs = ",len(pois) - pois.sum()
433
434
435 # Let's remove non usefull features
436 RemoveFeature(data_dict,'loan_advances')
437 RemoveFeature(data_dict,'director_fees')
438 RemoveFeature(data_dict,'restricted_stock_deferred')
439
440 RemovePositiveNegative(data_dict,'deferral_payments',removeNegative=True)
441 RemovePositiveNegative(data_dict,'restricted_stock',removeNegative=True)
442 RemovePositiveNegative(data_dict,'total_stock_value',removeNegative=True)
443
444 allFeatures = ['salary', 'bonus', 'long_term_incentive', 'deferred_income', 'deferral_payments', 'other', 'expenses', 'total_payments', 'ex
445 df = pandas.DataFrame(featureFormat(data_dict, allFeatures, sort_keys = False,remove_NaN=False),columns=allFeatures)
446
447 if showunivariate:
448     ShowHist(data_dict,allFeatures)
449
450 # Ok, now we do have out list of features and data dictionary (it includes POIs, so we need to remove one)
451 print "Number of features = ", len(allFeatures) -1
452
453 print "Correlation between features"
454 if showheatmap:
455     plt.figure(GetSequence())
456     sns.heatmap(df.corr(),xticklabels=allFeatures,yticklabels=allFeatures)
457     sns.plt.show()
458     print df.corr()
459
460 if showcorrelation:
461     ShowScatter(data_dict,"deferral_payments", "deferred_income")
462     ShowScatter(data_dict,"restricted_stock", "total_stock_value")
463
464     ShowScatter(data_dict,"to_messages", "from_this_person_to_poi")
465     ShowScatter(data_dict,"to_messages", "shared_receipt_with_poi")
466
467 # Let's add a new feature
468 CreateRatio(data_dict)
469
470 # Let's have a look on new ratio/poi possible correlation
471 df = pandas.DataFrame(featureFormat(data_dict, ['poi_to_ratio','poi_from_ratio','poi'], sort_keys = False,remove_NaN=False),columns=['po
472
473 if showheatmap:
474     plt.figure(GetSequence())
475     sns.heatmap(df.corr(),xticklabels=['poi_to_ratio','poi_from_ratio','poi'],yticklabels=['poi_to_ratio','poi_from_ratio','poi'])
476     sns.plt.show()
477     print df.corr()
478
479
480
481
482 ### Let's do the tuning
483 optim_dataset = ''
484 optim_log = ''
485 optim_val_strategy = ''
486 optim_algo = ''
487 optim_f1 = 0
488
489 if performalltunings:
490     for dataset in ['full','full_ratio','limited']:
491         for log in ['yes','no']:
492             for val_strategy in ['basic','kfold']:
493                 for algo in ['NB','SVM','DT']:
494                     if dataset == 'full':
495                         allFeatures = ['poi','salary', 'bonus', 'long_term_incentive', 'deferred_income', 'deferral_payments','other','e
496                     elif dataset == 'full_ratio':
497                         allFeatures = ['poi','salary', 'bonus', 'long_term_incentive', 'deferred_income', 'deferral_payments','other','e
498                     elif dataset == 'limited':
499                         allFeatures = ['poi','total_payments', 'total_stock_value', 'to_messages', 'from_poi_to_this_person', 'from_messa
500
501                     data = data_dict
502
503                     if log == 'yes':
504                         MakeLog(data)
505
506                     finaldata = featureFormat(data, allFeatures, sort_keys = False,remove_NaN=True)
507
508                     print "Test = dataset / log data / val strategy / algo = " , dataset, " / " , log, " / " , val_strategy, " / " ,algo
509                     if algo == "NB":
510                         clf, f1 = TuneGNB(data=finaldata,features_list=allFeatures,max_features=len(allFeatures),val_strategy=val_strategy)
511                     if algo == "DT":
512                         clf, f1 = TuneDT(data=finaldata,features_list=allFeatures,max_features=len(allFeatures),val_strategy=val_strategy)
513                     if algo == "SVM":
514                         clf, f1 = TuneSVM(data=finaldata,features_list=allFeatures,max_features=len(allFeatures),val_strategy=val_strategy)

```

```

515
516         if f1 > optim_f1:
517             optim_f1=f1
518             optim_algo = algo
519             optim_dataset = dataset
520             optim_log = log
521             optim_val_strategy = val_strategy
522             optim_clf = clf
523
524         print "F1 vs Optim F1",f1,optim_f1
525
526
527     print "Optim f1          = ",optim_f1
528     print "Optim algo       = ",optim_algo
529     print "Optim dataset    = ",optim_dataset
530     print "Optim log        = ",optim_log
531     print "Optim val strategy = ",optim_val_strategy
532
533     data = data_dict
534     if optim_dataset == 'full':
535         allFeatures = ['poi','salary', 'bonus', 'long_term_incentive', 'deferred_income', 'deferral_payments','other','expenses','total_
536     elif optim_dataset == 'full_ratio':
537         allFeatures = ['poi','salary', 'bonus', 'long_term_incentive', 'deferred_income', 'deferral_payments','other','expenses','total_
538     elif optim_dataset == 'limited':
539         allFeatures = ['poi','total_payments', 'total_stock_value','to_messages', 'from_poi_to_this_person', 'from_messages', 'from_this
540
541     if optim_log == 'yes':
542         MakeLog(data)
543
544     dump_classifier_and_data(optim_clf, data, allFeatures)
545
546
547
548
549     for rnd in [20,30,40,50,60]:
550         labels, features = targetFeatureSplit(featureFormat(data, allFeatures, sort_keys = False,remove_NaN=True))
551         features_train,features_test,labels_train,labels_test = train_test_split(features,labels,test_size=0.3,random_state=rnd)
552
553         prediction = optim_clf.predict(features_test)
554         print classification_report(labels_test,prediction)
555
556

```

Have a question about your review? Email us at [review-support@udacity.com](mailto:review-support@udacity.com) and include the link to this review.

RETURN TO PATH

[Student FAQ](#)