



PROJECT

Identify Fraud from Enron Email

A part of the Data Analyst Nanodegree Program

PROJECT REVIEW

CODE REVIEW 3

NOTES

▼ poi_id.py 3

```

1  #!/usr/bin/python
2  from __future__ import division
3  from matplotlib import pyplot as plt
4  from time import time
5  import numpy as np
6  import matplotlib
7  import pandas as pd
8  import seaborn as sns
9  from sklearn.metrics import recall_score
10 from sklearn.metrics import precision_score
11 from sklearn import svm
12 from sklearn import neighbors
13 from sklearn.svm import SVC
14 from sklearn.ensemble import RandomForestClassifier
15 from sklearn.model_selection import train_test_split
16 from sklearn.model_selection import GridSearchCV
17 from sklearn.naive_bayes import GaussianNB
18 import collections
19
20 from sklearn.preprocessing import StandardScaler, RobustScaler, MinMaxScaler
21 from sklearn.feature_selection import SelectKBest, f_classif
22
23 import sys
24 import pickle
25 sys.path.append("../tools/")
26
27 from feature_format import featureFormat, targetFeatureSplit
28 from tester import *
29
30 ### Task 1: Select what features you'll use.
31 def load_data():
32
33     ### Load the dictionary containing the dataset
34     data_dict = {}
35     with open("final_project_dataset.pkl", "r") as data_file:
36         data_dict = pickle.load(data_file)
37
38     return data_dict
39
40
41 ### Task 2: Remove outliers
42 def remove_outlier(data_dict):
43     data_dict.pop( "TOTAL", 0 )
44

```

SUGGESTION

I will leave hints on a few other outliers we could remove.

```

45     return data_dict
46
47     ## Task 3: Create new feature(s)
48     ## Store to my_dataset for easy export below.
49     def create_new_features(data_dict, features_list):
50

```



SUGGESTION

I have two suggestions here:

1. Further comment each new method that has been created using [PEP 257](#), and
2. Move the auxiliary methods to a separate file, thus making the main code shorter and easier to read.

```

51     bProc_shared_receipt_with_poi = False
52     bProc_ratio_stock_option = False
53     bProc_salary_bonus_ratio = False
54     for person, item in data_dict.items():
55
56         ## Checking if there's any correlation if the person sends more email to poi vs all the messages to indicate i
57         key = 'ratio_from_poi_fr_msg'
58         if not(key in item) and not(key in features_list):
59             b_ratio_from_poi_fr_msg = True
60             features_list.append(key)
61
62         if b_ratio_from_poi_fr_msg:
63             item[key] = computeFraction(item['from_poi_to_this_person'], item['from_messages'])
64
65         ## Also checking the other way around
66         key = 'ratio_to_poi_to_msg'
67         if not(key in item) and not(key in features_list):
68             b_ratio_to_poi_to_msg = True
69             features_list.append(key)
70
71         if b_ratio_to_poi_to_msg:
72             item[key] = computeFraction(item['from_this_person_to_poi'], item['to_messages'])
73
74         ## and then check if total combined communication with the poi
75         key = 'total_to_poi_from_poi'
76         if not(key in item) and not(key in features_list):
77             b_total_to_poi_from_poi = True
78             features_list.append(key)
79
80         if b_total_to_poi_from_poi:
81             item[key] = computeAddition(item['ratio_to_poi_to_msg'], item['ratio_from_poi_fr_msg'])
82
83         ## Seeing a high correlation with stock and salary hence using both as a new total feature
84         key = 'total_salary_stock'
85         if not(key in item) and not(key in features_list):
86             b_total_salary_stock = True
87             features_list.append(key)
88
89         if (b_total_salary_stock):
90             item[key] = computeAddition(item['salary'], item['total_stock_value'])
91
92
93         ## Create a new feature which combine both bonus and stock value
94         key = 'total_bonus_stock'
95         if not(key in item) and not(key in features_list):
96             b_total_bonus_stock = True
97             features_list.append(key)
98
99         if b_total_bonus_stock:
100             item[key] = computeAddition(item['total_stock_value'], item['bonus'])
101
102     return data_dict, features_list
103
104
105
106
107     ## Task 4: Try a variety of classifiers
108     ## Please name your classifier clf for easy export below.
109     ## Note that if you want to do PCA or other multi-stage operations,
110     ## you'll need to use Pipelines. For more info:

```

```

111 ### http://scikit-learn.org/stable/modules/pipeline.html
112 def apply_classifiers():
113     clf = svm.LinearSVC()
114     clf2 = neighbors.KNeighborsClassifier(n_neighbors=4, weights="distance", leaf_size=30, algorithm='brute')
115     clf3 = GaussianNB()
116     clf4 = RandomForestClassifier(n_estimators=100)
117
118     clfs = collections.OrderedDict()
119     clfs['LinearSVC'] = clf
120     clfs['KNearestNeighbour'] = clf2
121     clfs['NaiveBayes'] = clf3
122     clfs['RandomForest'] = clf4
123
124     return clfs
125
126
127 ### Task 5: Tune your classifier to achieve better than .3 precision and recall
128 using our testing script. Check the tester.py script in the final project
129 folder for details on the evaluation method, especially the test\_classifier
130 function. Because of the small size of the dataset, the script uses
131 stratified shuffle split cross validation. For more info:
132 ### http://scikit-learn.org/stable/modules/generated/sklearn.cross\_validation.StratifiedShuffleSplit.html
133 def tune_classifier(features_train, labels_train):
134     # Example starting point. Try investigating other evaluation techniques!
135     #features_train, features_test, labels_train, labels_test = \
136     #     train_test_split(features, labels, test_size=0.3, random_state=42)
137
138     tuned_clfs = collections.OrderedDict()
139     svrs = collections.OrderedDict()
140
141     Cs = [0.01, 0.1, 1, 10, 100, 1000, 10000]
142     param_grid = {'C': Cs}
143     svr = svm.LinearSVC()
144     svrs['Tuned LinearSVC'] = [svr, param_grid]
145
146     n_neighbours = range(1,10)
147     weights = ['distance', 'uniform']
148     leaf_sizes = range(1,10)
149     algorithms = ['kd_tree', 'brute', 'ball_tree']
150
151     param_grid = {'n_neighbors': n_neighbours, 'weights': weights, 'algorithm': algorithms, 'leaf_size': leaf_sizes}
152     svr = neighbors.KNeighborsClassifier()
153     svrs['Tuned KNearestNeighbour'] = [svr, param_grid]
154
155     n_estimators = [1,5,10,100]
156     max_features = ['sqrt', 'log2']
157     min_samples_splits = range(2,10)
158     min_samples_leafs = range(1,10)
159
160     param_grid = {'n_estimators': n_estimators, 'max_features': max_features, 'min_samples_split': min_samples_splits, '
161     svr = RandomForestClassifier()
162     svrs['Tuned RandomForest'] = [svr, param_grid]
163
164     for key, item in svrs.iteritems():
165         clf_gs = GridSearchCV(item[0], item[1])
166         clf_gs = clf_gs.fit(features_train, labels_train)
167         clf_gs.best_params_
168         print "Best estimator found by grid search:"
169         print clf_gs.best_estimator_, "\n"
170
171         tuned_clfs[key] = clf_gs.best_estimator_
172
173     return tuned_clfs
174
175 ### Task 6: Dump your classifier, dataset, and features\_list so anyone can
176 check your results. You do not need to change anything below, but make sure
177 that the version of poi\_id.py that you submit can be run on its own and
178 generates the necessary .pkl files for validating your results.
179 def dump_classifier_results(clf, dataset, features_list):
180     dump_classifier_and_data(clf, dataset, features_list)
181
182
183 ### Common functions
184 def plotScatter(data, xlab, ylab, xidx=0, yidx=1):
185     fig = plt.figure()
186     ax = fig.add_subplot(111)
187     #ax.set_xticks(np.arange(0, 2, 1))
188
189
190     ### your code below
191     for point in data:

```

```

192     x = point[xidx]
193     y = point[yidx]
194     if point[0] == 1.0:
195         selcolor = 'r'
196     else:
197         selcolor = 'b'
198     plt.scatter( x, y, color=selcolor, alpha=.4 )
199
200     plt.xlabel(xlab)
201     plt.ylabel(ylab)
202     plt.show()
203
204 def show_data_overview(data_dict, features_list):
205     print "\ntotal number of data points: ", len(data_dict)
206     print "allocation across classes (POI/non-POI): ", sum(x['poi'] for x in data_dict.values()), "/", sum(x['poi']==0
207     print "no of features per person: ", len(data_dict[data_dict.keys()[0]])
208     print "number of features used: ", len(features_list)
209     print "selected features:", " ", ".join(features_list)
210
211 def select_features(data_dict, filter_pct):
212     features_list = []
213     print "Are there features with many missing values? etc."
214     sorted_nan_dict = {}
215     for item in data_dict[data_dict.keys()[0]]:
216         sorted_nan_dict[item] = sum(x[item]=="NaN" for x in data_dict.values())
217         #print " ", item, ":", sum(x[item]=="NaN" for x in data_dict.values())
218
219     # remove any features that have missing values above the given threshold
220     for k,v in sorted(sorted_nan_dict.iteritems(), key=lambda(k,v) : (v,k), reverse = True):
221         pct = (v / len(data_dict))
222         print " ", k, ":", v, "Missing data percentage:", round(pct*100,2), "%"
223         # filter out email address which is not a numeric value
224         if(pct <= filter_pct and k != "email_address"):
225             features_list.append(k)
226
227     return features_list
228
229 def computeFraction( poi_messages, all_messages ):
230     """ given a number messages to/from POI (numerator)
231     and number of all messages to/from a person (denominator),
232     return the fraction of messages to/from that person
233     that are from/to a POI
234     """
235
236
237     ### you fill in this code, so that it returns either
238     ### the fraction of all messages to this person that come from POIs
239     ### or
240     ### the fraction of all messages from this person that are sent to POIs
241     ### the same code can be used to compute either quantity
242
243     ### beware of "NaN" when there is no known email address (and so
244     ### no filled email features), and integer division!
245     ### in case of poi_messages or all_messages having "NaN" value, return 0.
246     fraction = 0.
247     if (poi_messages != "NaN" and all_messages != "NaN"):
248         fraction = poi_messages / all_messages
249
250     return fraction
251
252 def computeAddition( var1, var2):
253     total = 0
254     if var1 != "NaN" and var2 != "NaN":
255         total = var1 + var2
256
257     return total
258
259 def plotCorrMatrix(seldata, features_list):
260     sns.set(style="white")
261
262     d = pd.DataFrame(data=seldata, columns=features_list)
263     #print d
264
265     # Compute the correlation matrix
266     corr = d.corr()
267
268     # Generate a mask for the upper triangle
269     mask = np.zeros_like(corr, dtype=np.bool)
270     mask[np.triu_indices_from(mask)] = True
271
272     # Set up the matplotlib figure

```

```

273 f, ax = plt.subplots(figsize=(11, 9))
274
275 # Generate a custom diverging colormap
276 cmap = sns.diverging_palette(220, 10, as_cmap=True)
277
278 # Draw the heatmap with the mask and correct aspect ratio
279 sns.heatmap(corr, mask=mask, cmap=cmap, vmax=1., center=0,
280             square=True, linewidths=.5, cbar_kws={"shrink": .5});
281
282 plt.show()
283
284 def split_to_label_features(data_dict, features_list):
285     """ Extract features and labels from dataset for local testing
286     data = featureFormat(data_dict, features_list, sort_keys = True)
287     labels, features = targetFeatureSplit(data)
288
289     return labels, features
290
291 def autoselect_features(data_dict, features_list):
292     labels, features = split_to_label_features(data_dict, features_list)
293
294     # looking for the top 5 features to use
295     selector = SelectKBest(f_classif, k = 5)
296     selector.fit(features, labels)
297     scores = zip(features_list[1:], selector.scores_)
298     print 'SelectKBest scores: ', sorted(scores, key=lambda x: x[1], reverse=True)
299
300     #print len(features)
301     final_features = selector.transform(features)
302     #print len(final_features)
303     return final_features
304
305 def apply_robust_feature_scaling(data):
306     scaler = RobustScaler()
307     scaler.fit(data)
308     return scaler.transform(data)
309
310 def apply_standard_feature_scaling(data):
311     scaler = StandardScaler()
312     scaler.fit(data)
313     return scaler.transform(data)
314
315 def apply_minmax_feature_scaling(data):
316     scaler = MinMaxScaler()
317     scaler.fit(data)
318     return scaler.transform(data)
319
320
321 def clf_score_and_evaluate(clf, features_train, labels_train, features_test, labels_test):
322     clf.fit(features_train, labels_train)
323     pred = clf.predict(features_test)
324     sel_avg = 'macro'
325     display_precision = 5
326     accuracy = round(clf.score(features_test, labels_test), display_precision)
327     #print "accuracy: ", round(accuracy, 3)
328
329     precision = round(precision_score(labels_test, pred, average=sel_avg), display_precision)
330     #print "precision: ", round(precision, 3)
331
332     recall = round(recall_score(labels_test, pred, average=sel_avg), display_precision)
333     #print "recall: ", round(recall, 3)
334
335     return accuracy, precision, recall
336
337 def compare_algorithms(clfs, features_train, labels_train, features_test, labels_test, sc_features_train, sc_labels_train):
338
339     datamatrix = []
340     for name, clf in clfs.iteritems():
341
342         clf.fit(features_train, labels_train)
343         accuracy, precision, recall = clf_score_and_evaluate(clf, features_train, labels_train, features_test, labels_test)
344
345         #print "\nRescaled features:"
346         clf.fit(sc_features_train, sc_labels_train)
347         sc_accuracy, sc_precision, sc_recall = clf_score_and_evaluate(clf, sc_features_train, sc_labels_train, sc_features_test, sc_labels_test)
348
349         datarow = [name, accuracy, precision, recall, sc_accuracy, sc_precision, sc_recall]
350         datamatrix.append(datarow)
351
352     return datamatrix

```

```
353 def test_algorithms(clfs, dataset, features_list):
354     for name, clf in clfs.iteritems():
355         test_classifier(clf, dataset, features_list)
356
357
358 def apply_cross_validation(features, labels, scaled_features):
359     features_train, features_test, labels_train, labels_test = train_test_split(features, labels, test_size=0.3, random_state=42)
```

SUGGESTION

I will also leave a suggestion about the validation scheme on the other tab.

```
360
361     sc_features_train, sc_features_test, sc_labels_train, sc_labels_test = train_test_split(scaled_features, labels, test_size=0.3, random_state=42)
362
363     return features_train, features_test, labels_train, labels_test, sc_features_train, sc_features_test, sc_labels_train, sc_labels_test
```

► readMe.rmd

RETURN TO PATH

Rate this review

[Student FAQ](#)