# UDACITY

PROJECT

## Identify Fraud from Enron Email

A part of the Data Analyst Nanodegree Program

| PROJECT REVIEW |
| --- |
| CODE REVIEW  7 |
| NOTES |

▼ poi_id.py        7

```python
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 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 select_features(data_dict, filter_pct):
32      """ """
33      features_list = []
34      print "Are there features with many missing values? etc."
35      sorted_nan_dict = {}
36      for item in data_dict[data_dict.keys()[0]]:
37          sorted_nan_dict[item] = sum(x[item]=="NaN" for x in data_dict.values())
38
39      # remove any features that have missing values above the given threshold
40      for k,v in sorted(sorted_nan_dict.iteritems(), key=lambda(k,v) : (v,k), \
41          reverse = True):
42          pct = (v / len(data_dict))
43          print " ", k, ": ", v, ",Missing data percentage:", round(pct*100,2), "%"
44          # filter out email address which is not a numeric value
45          if(pct <= filter_pct and k != "email_address"):
46              features_list.append(k)
47
        #rearrange the list order to have poi in the first position of the list
```

```
49      features_list =  [features_list[-1]]+features_list[0:len(features_list)-1:1]
50
51      return features_list
52
53
54  ### Task 2: Remove outliers
55  def remove_outlier(data_dict):
```

▲

AWESOME

Excellent!

```
56      data_dict.pop( "TOTAL", 0 )
57      data_dict.pop("THE TRAVEL AGENCY IN THE PARK", 0)
58      data_dict.pop("LOCKHART EUGENE E", 0)
59
60      return data_dict
61
62  ### Task 3: Create new feature(s)
63  def create_new_features(data_dict, features_list):
```

▲

AWESOME

Nice work engineering your new features!

```
64
65      b_ratio_from_poi_fr_msg = False
66      b_ratio_to_poi_to_msg = False
67      b_total_to_poi_from_poi = False
68      b_total_salary_stock = False
69      b_total_bonus_stock = False
70      for person, item in data_dict.items():
71
72          ## Checking if there's any correlation if the person sends more email
73          ## to poi vs all the messages to indicate if he can also be poi
74          key = 'ratio_from_poi_fr_msg'
75          if not(key in item) and not(key in features_list):
76              b_ratio_from_poi_fr_msg = True
77              features_list.append(key)
78
79          if b_ratio_from_poi_fr_msg:
80              item[key] = computeFraction(item['from_poi_to_this_person'],
81                  item['from_messages'])
82
83          ## Also checking the other way around
84          key = 'ratio_to_poi_to_msg'
85          if not(key in item) and not(key in features_list):
86              b_ratio_to_poi_to_msg = True
87              features_list.append(key)
88
89          if b_ratio_to_poi_to_msg:
90              item[key] = computeFraction(item['from_this_person_to_poi'],
91                  item['to_messages'])
92
93          ## and then check if total combined communication with the poi
94          key = 'total_to_poi_from_poi'
95          if not(key in item) and not(key in features_list):
96              b_total_to_poi_from_poi = True
97              features_list.append(key)
98
99          if b_total_to_poi_from_poi:
100             item[key] = computeAddition(item['ratio_to_poi_to_msg'],
101                 item['ratio_from_poi_fr_msg'])
102
103         ## Seeing a high correlation with stock and salary hence using both as
104         ## a new total feature
105         key = 'total_salary_stock'
106         if not(key in item) and not(key in features_list):
107             b_total_salary_stock = True
108             features_list.append(key)
109
110         if (b_total_salary_stock):
111             item[key] = computeAddition(item['salary'],
112                 item['total_stock_value'])
113
114
```

```
115          ## Create a new feature which combine both bonus and stock value
116          key = 'total_bonus_stock'
117          if not(key in item) and not(key in features_list):
118              b_total_bonus_stock = True
119              features_list.append(key)
120
121          if b_total_bonus_stock:
122              item[key] = computeAddition(item['total_stock_value'],item['bonus'])
123
124
125      return data_dict, features_list
126
127
128
129  ### Task 4: Try a variety of classifiers
130  def apply_classifiers():
```

AWESOME

Well done attempting so many classifiers

```
131      clf = svm.LinearSVC()
132      clf2 = neighbors.KNeighborsClassifier(n_neighbors=4, weights="distance", \
133          leaf_size=30, algorithm='brute')
134      clf3 = GaussianNB()
135      clf4 = RandomForestClassifier(n_estimators=100)
136
137      clfs = collections.OrderedDict()
138      clfs['LinearSVC'] = clf
139      clfs['KNearestNeighbour'] = clf2
140      clfs['NaiveBayes'] = clf3
141      clfs['RandomForest'] = clf4
142
143      return clfs
144
145
146  ### Task 5: Tune your classifier to achieve better than .3 precision and recall
147  ### using our testing script.
148  def tune_classifier(features_train, labels_train):
149
150      tuned_clfs = collections.OrderedDict()
151      svrs = collections.OrderedDict()
152
153      Cs = [0.01, 0.1, 1, 10, 100, 1000, 10000]
154      param_grid = {'C': Cs}
155      svr = svm.LinearSVC()
156      svrs['Tuned LinearSVC'] = [svr, param_grid]
157
158      n_neighbours = range(1,10)
159      weights = ['distance', 'uniform']
160      algorithms = ['kd_tree', 'brute', 'ball_tree']
161
162      param_grid = {'n_neighbors': n_neighbours, 'algorithm':algorithms, \
163                  'weights':weights}
164      svr = neighbors.KNeighborsClassifier()
165      svrs['Tuned KNearestNeighbour'] = [svr, param_grid]
166
167      n_estimators = range(10,100,10)
168
169      param_grid = {'n_estimators': n_estimators}
170      svr = RandomForestClassifier()
171      svrs['Tuned RandomForest'] = [svr, param_grid]
172
173      for key, item in svrs.iteritems():
174          clf_gs = GridSearchCV(item[0], item[1], scoring='recall_macro')
```

AWESOME

Good optimization of GridCV!

```
175          clf_gs = clf_gs.fit(features_train, labels_train)
176          clf_gs.best_params_
177          print "Best estimator found by grid search:"
178          print clf_gs.best_estimator_, "\n"
179
180          tuned_clfs[key] = clf_gs.best_estimator_
```

```python
181
182    return tuned_clfs
183
184 ### Task 6: Dump your classifier, dataset, and features_list so anyone can
185 ### check your results.
186 def dump_classifier_results(clf, dataset, features_list):
187     dump_classifier_and_data(clf, dataset, features_list)
188
189
190 def load_data():
191
192     """ Load the dictionary containing the dataset """
193     data_dict = {}
194     with open("final_project_dataset.pkl", "r") as data_file:
195         data_dict = pickle.load(data_file)
196
197     return data_dict
198
199 def plotScatter(data, xlab, ylab, xidx=0, yidx=1):
200     """ plot a scatter plot """
201     fig = plt.figure()
202     ax = fig.add_subplot(111)
203
204     for point in data:
205         x = point[xidx]
206         y = point[yidx]
207         if point[0] == 1.0:
208             selcolor = 'r'
209         else:
210             selcolor = 'b'
211         plt.scatter( x, y, color=selcolor, alpha=.4 )
212
213     plt.xlabel(xlab)
214     plt.ylabel(ylab)
215     plt.show()
216
217 def show_data_overview(data_dict, features_list):
218     """ show the overall data structure """
219     print "\ntotal number of data points: ", len(data_dict)
220     print "allocation across classes (POI/non-POI): ", \
221         sum(x['poi'] for x in data_dict.values()), "/", \
222         sum(x['poi']==0 for x in data_dict.values())
223     print "no of features per person: ", len(data_dict[data_dict.keys()[0]])
224     print "number of features used: ", len(features_list)
225     print "selected features:", ", ".join(features_list)
226
227
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     fraction = 0.
236     if (poi_messages != "NaN" and all_messages != "NaN"):
237         fraction = poi_messages / all_messages
238
239     return fraction
240
241 def computeAddition( var1, var2):
242     """ does the addition given two list and exclude NaN values """
243     total = 0
244     if var1 != "NaN" and var2 != "NaN":
245         total = var1 + var2
246
247     return total
248
249 def plotCorrMatrix(seldata, features_list):
```

AWESOME

Great work with this visualization!

```python
250     """ plot correlation matrix """
251     sns.set(style="white")
252
253     d = pd.DataFrame(data=seldata, columns=features_list)
```

```
254        #print d
255
256        # Compute the correlation matrix
257        corr = d.corr()
258
259        # Generate a mask for the upper triangle
260        mask = np.zeros_like(corr, dtype=np.bool)
261        mask[np.triu_indices_from(mask)] = True
262
263        # Set up the matplotlib figure
264        f, ax = plt.subplots(figsize=(14, 12))
265
266        # Generate a custom diverging colormap
267        cmap = sns.diverging_palette(220, 10, as_cmap=True)
268
269        # Draw the heatmap with the mask and correct aspect ratio
270        sns.heatmap(corr, mask=mask, cmap=cmap, vmax=1., center=0, \
271                square=True, linewidths=.5, cbar_kws={"shrink": .5}, annot=True);
272
273        plt.show()
274
275 def split_to_label_features(data_dict, features_list):
276        """ Extract features and labels from dataset for local testing """
277        data = featureFormat(data_dict, features_list, sort_keys = False)
278        labels, features = targetFeatureSplit(data)
279
280        return labels, features
281
282 def autoselect_features(data_dict, features_list, sel_k = 5, bShow=True):
283        """
284        select features based on anova f stats and can take k as no of features
285        to be returned.
286        """
287        labels, features = split_to_label_features(data_dict, features_list)
288
289        # looking for the top k features to use
290        selector = SelectKBest(f_classif, k = sel_k)
291        selector.fit(features, labels)
292
293        # print the f-stat score sorted in descending order
294        if bShow:
295            scores = zip(features_list[1:],[round(elem,4) for elem in selector.scores_])
```

AWESOME

Nice work getting the scores!

```
296            print 'SelectKBest scores: ', sorted(scores, key=lambda x: x[1], \
297                reverse=True)
298
299        final_features = selector.transform(features)
300
301        # Get idxs of columns to keep
302        idxs_selected = selector.get_support(indices=True)
303        selected_features_list = [features_list[i+1] for i in idxs_selected]
304
305        return final_features, labels, selected_features_list
306
307
308 def apply_robust_feature_scaling(data):
309        scaler = RobustScaler()
310        scaler.fit(data)
311        return scaler.transform(data)
312
313 def apply_standard_feature_scaling(data):
314        scaler = StandardScaler()
315        scaler.fit(data)
316        return scaler.transform(data)
317
318 def apply_minmax_feature_scaling(data):
319        scaler = MinMaxScaler()
320        scaler.fit(data)
321        return scaler.transform(data)
322
323
324 def clf_score_and_evaluate(cur_clf,features_train, labels_train,features_test, \
```

SUGGESTION

It is great you validate your algorithm, however, there are ways to improve your validation process.

An important part of the validation process is cross-validation, cross-validation generates pairs of train/test sets and these pairs are used to fit/pr
many pairs of train/test sets as required in order to generalize results over the entire dataset.

In your validation process implemented it is just generated **one** train/test pair used to fit/predict/evaluate your algorithm. As you can see, this ma

STEPS TO ENHANCE YOUR VALIDATION PROCESS:

1. Create local variables used to store your score values generated in each split, for example:

```
score_all = []
precision_all = []
recall_all = []
```

2. Include in the for loop the fit/prediction/evaluation process and append results to local variables:

```
kf = #Cross-Validation Object
for train_indices, test_indices in kf:
    #make training and testing sets
    features_train= [features[ii] for ii in train_indices]
    features_test= [features[ii] for ii in test_indices]
    labels_train=[labels[ii] for ii in train_indices]
    labels_test=[labels[ii] for ii in test_indices]
    clf = #Algorithm to validate
    clf.fit(features_train,labels_train)
    pred = clf.predict(features_test)
    score_all.append(clf.score(features_test,labels_test))
    precision_all.append(precision_score(labfrom sklearn.preprocessing import MinMaxScalerels_test,pred))
    recall_all.append(recall_score(labels_test,pred))
```

3. Once the for-loop ends, calculate your validation metrics from your local variables:

```
precision = numpy.average(precision_all)
recall = numpy.average(recall_all)
score = numpy.average(score_all)
```

**Important:** Note this dataset is unbalanced and small, so the splitting technique more appropriated is StratifiedShuffleSplit, since it creates splits
complete set. Have a look at this link for further information.

**Important 2:** Note cross-validation is an iterative process where train/test are generated until results converge, that is, are not dependent on the
process is repeated 1000 times!.

```
325    labels_test):
326    cur_clf.fit(features_train, labels_train)
327    pred = cur_clf.predict(features_test)
328    sel_avg = 'macro'
329    display_precision = 5
330    accuracy = round(cur_clf.score(features_test,labels_test),display_precision)
331    #print "accuracy: ", round(accuracy,3)
332
333    precision = round(precision_score(labels_test, pred, average=sel_avg), \
334        display_precision)
335    #print "precision: ", round(precision, 3)
336
337    recall = round(recall_score(labels_test, pred, average=sel_avg), \
338        display_precision)
339    #print "recall: ", round(recall, 3)
340
341    return accuracy, precision, recall
342
343 def compare_algorithms(clfs, features_train, labels_train, features_test, \
344    labels_test, sc_features_train, sc_labels_train,sc_features_test, \
345    sc_labels_test):
346
347    datamatrix = []
348    for name, cur_clf in clfs.iteritems():
349
350        accuracy, precision, recall = clf_score_and_evaluate(cur_clf, \
351            features_train, labels_train,features_test, labels_test)
352
353        #print "\nRescaled features:"
354        sc_accuracy, sc_precision, sc_recall = clf_score_and_evaluate(cur_clf, \
355            sc_features_train, sc_labels_train,sc_features_test, sc_labels_test)
356
357        datarow = [name,accuracy, precision, recall,sc_accuracy, sc_precision,\
358            sc_recall]
```

```
359            datamatrix.append(datarow)
360
361        return datamatrix
362
363  def test_algorithms(clfs, dataset, features_list):
364      for name, clf in clfs.iteritems():
365          test_classifier(clf, dataset, features_list)
366
367  def apply_cross_validation(features, labels, scaled_features):
368      features_train, features_test, labels_train, labels_test = \
369      train_test_split(features, labels, test_size=0.3, random_state=42)
370
371      sc_features_train, sc_features_test, sc_labels_train, sc_labels_test = \
372      train_test_split(scaled_features, labels, test_size=0.3, random_state=42)
373
374      return features_train, features_test, labels_train, labels_test, \
375      sc_features_train, sc_features_test, sc_labels_train, sc_labels_test
376
377  def get_features_score(features, labels):
378      """ use linear svc to compare the features score """
379      features_train, features_test, labels_train, labels_test = \
380          train_test_split(features, labels, test_size=0.3, random_state=42)
381
382      clf = neighbors.KNeighborsClassifier(n_neighbors=4, weights="distance", \
383          leaf_size=30, algorithm='brute')
384
385      accuracy, precision, recall = clf_score_and_evaluate(clf, \
386          features_train, labels_train,features_test, labels_test)
387
388      return accuracy, precision, recall
```

▶ readMe.rmd

RETURN TO PATH

Rate this review

Student FAQ