



[Return to "Data Analyst Nanodegree" in the classroom](#)

Identify Fraud from Enron Email

REVIEW

CODE REVIEW 4

HISTORY

▼ poi_id.py 4

```
1 #!/usr/bin/env python
2 # coding: utf-8
3
4 # ## Projeto Final - DS2
5
6 # ### Import
7
8 # In[232]:
9
10
11 import sys
12 import pickle
13 import pandas as pd
14 import numpy as np
15 import warnings
16 import matplotlib.pyplot as plt
17 sys.path.append("../tools/")
18 warnings.filterwarnings("ignore")
19
20
21 # In[233]:
22
23
24 from sklearn.pipeline import Pipeline
25
26 from sklearn import preprocessing
27 from sklearn.decomposition import PCA
28 from sklearn.model_selection import GridSearchCV, train_test_split,
29 from sklearn.tree import DecisionTreeClassifier
```

```

29 from sklearn.feature_selection import SelectKBest
30 from sklearn.externals import joblib
31 from sklearn.linear_model import LogisticRegression
32 from sklearn.naive_bayes import GaussianNB
33 from sklearn.grid_search import GridSearchCV
34 from sklearn.cross_validation import StratifiedShuffleSplit, train_test_split
35 from sklearn.neighbors import KNeighborsClassifier
36 from sklearn.svm import SVC
37 from sklearn.metrics import accuracy_score
38 from sklearn.preprocessing import MinMaxScaler
39 from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
40
41
42 # In[234]:
43
44
45 from time import time
46 from feature_format import featureFormat, targetFeatureSplit
47 from tester import dump_classifier_and_data, test_classifier
48 scaler = MinMaxScaler()
49
50
51 # ### Funções
52
53 # In[307]:
54
55
56 def monta_grafico(feet_x, feat_y, titulo, dicionario, cor):

```

SUGGESTION

Bom trabalho ao criar funções auxiliares! Minha sugestão é de movê-las para arquivos auxiliares.

```

57     # Criar um grafico scatter das fetures passadas no parametro
58     features = ['poi', feat_x, feat_y]
59     data = featureFormat(dicionario, features)
60
61     plt.figure(figsize=(16,7))
62
63     for point in data:
64         x = point[1]
65         y = point[2]
66         if point[0]:
67             if cor == 1:
68                 plt.scatter(x, y, color="red", marker="*")
69             else:
70                 plt.scatter(x, y, color="green", marker=".")
71         else:
72             if cor == 1:
73                 plt.scatter(x, y, color='blue', marker=".")
74             else:
75                 plt.scatter(x, y, color="orange", marker="*")
76
77     plt.title(titulo, fontsize=20)
78     plt.xlabel(feet_x, fontsize=18)
79     plt.ylabel(feat_y, fontsize=18)
80     pic = feat_x + feat_y + '.png'
81     plt.savefig(pic, transparent=True)

```

```

82     plt.show()
83     plt.show()
84
85
86 def monta_feature(features_list):
87     features_list = ['poi',
88                     'salary',
89                     'deferral_payments',
90                     'total_payments',
91                     'loan_advances',
92                     'bonus',
93                     'restricted_stock_deferred',
94                     'deferred_income',
95                     'total_stock_value',
96                     'expenses',
97                     'exercised_stock_options',
98                     'other',
99                     'long_term_incentive',
100                    'restricted_stock',
101                    'director_fees',
102                    'to_messages',
103                    'from_poi_to_this_person',
104                    'from_messages',
105                    'from_this_person_to_poi',
106                    'shared_receipt_with_poi']
107     return features_list
108
109 def nova_feature(dataset, features_list):
110     nova_feature = ["fraction_from_poi_email", "fraction_to_poi_email"]
111     num_features = ["from_poi_to_this_person", "from_this_person_to_poi"]
112     den_features = ["to_messages", "from_messages"]
113
114     for x in dataset:
115         data = dataset[x]
116
117         for i, feature in enumerate(nova_feature):
118             if data["poi"]:
119                 data[feature] = 'NaN'
120             else:
121                 message_poi = data[num_features[i]]
122                 messages_all = data[den_features[i]]
123                 fracao_messages = calcula_fracao(message_poi, messages_all)
124                 data[feature] = fracao_messages
125
126     return features_list + nova_feature
127
128
129 def testa_nova_feature(dataset, x, nova_feature):
130     num_features = ["from_poi_to_this_person", "from_this_person_to_poi"]
131     den_features = ["to_messages", "from_messages"]
132
133     print x, "\n- {} = {:.4f} ({} / {})\n- {} = {:.4f} ({} / {})\n".format(nova_feature[0],
134                                     data[num_features[0]] / data[den_features[0]],
135                                     nova_feature[1],
136                                     data[num_features[1]] / data[den_features[1]])
137
138
139 def pipeline_classificador(tipo, kbest, f_list):
140     # Contruir um pipeline e tune parameters via GridSearchCV
141
142     data = featureFormat(my_dataset, f_list, sort_keys=True)
143     labels, features = targetFeatureSplit(data)
144
145     # Usando o stratified shuffle split cross validation devido ao tamanho dos conjuntos
146     stratified_split_cross_validation = StratifiedShuffleSplit(labels, 500, test_size=0.2)

```

```

144
145 # Build pipeline
146 kbest = SelectKBest(k=kbest)
147 scaler = MinMaxScaler()
148 classifier = escolher_classificador(tipo)
149 pipeline = Pipeline(steps=[('minmax_scaler', scaler), ('feature_selection', kbest)
150
151 # Set parameters for random forest
152 parameters = []
153 if tipo == 'randomforest':
154     parameters = dict(randomforest__n_estimators=[25, 50],
155                       randomforest__min_samples_split=[2, 3, 4],
156                       randomforest__criterion=['gini', 'entropy'])
157 if tipo == 'logistic_regression':
158     parameters = dict(logistic_regression__class_weight=['balanced'],
159                       logistic_regression__solver=['liblinear'],
160                       logistic_regression__C=range(1, 5),
161                       logistic_regression__random_state=42)
162 if tipo == 'decisiontree':
163     parameters = dict(decisiontree__min_samples_leaf=range(1, 5),

```

SUGGESTION

Aqui seria interessante também tunar o parâmetro `k` do `feature_selection`!

```

164         decisiontree__max_depth=range(1, 5),
165         decisiontree__class_weight=['balanced'],
166         decisiontree__criterion=['gini', 'entropy'])
167
168 # Get optimized parameters for F1-scoring metrics
169 cv = GridSearchCV(pipeline, param_grid=parameters, scoring='f1', cv=stratified_sp:

```

AWESOME

Ótima implementação!

```

170 t0 = time()
171 cv.fit(features, labels)
172 print 'Tuning Classifier: %r' % round(time() - t0, 3)
173
174 return cv
175
176 def melhor_classificador(aux):
177     # Função para escolher um tipo de classificador
178     return {
179         'randomforest': RandomForestClassifier(),
180         'decisiontree': DecisionTreeClassifier(),
181         'logistic_regression': LogisticRegression(),
182         'gaussiannb': GaussianNB()
183     }.get(aux)
184
185 def calcula_fracao(message_poi, messages_all):
186     calc_fracao = 0.
187     if message_poi != "NaN" and messages_all != "NaN":
188         calc_fracao = float(message_poi) / messages_all
189     return calc_fracao

```

```

190
191 def sumariza_valores(dataset):
192     df_list = []
193     for key, y in dataset.items():
194         df_list.append(y)
195
196     df = pd.DataFrame(df_list, columns = dataset.items()[0][1].keys())
197
198     for i in df.columns:
199         df[i][df[i].apply(lambda i: True if str(i) == "NaN" else False)]=None
200
201     df = df.convert_objects(convert_numeric=True)
202     df.info()
203
204 def conta_poi(dataset):
205     poi_count = 0
206     for key, value in dataset.items():
207         if value['poi']:
208             poi_count += 1
209     return poi_count
210
211 def accuracy(new_features, features_list):
212     #Feature List
213     features_list = monta_feature(features_list)
214     if new_features == False:
215         print features_list
216
217     else:
218         features_list = nova_feature(data_dict_woo, features_list)
219         print ""
220         print features_list
221         print "Testando novos features adicionados:\n"
222         testa_nova_feature(data_dict_woo, "DIETRICH JANET R", features_list[-2:])
223
224     # Extrair as features e os labels do conjunto de dados
225     data = featureFormat(my_dataset, features_list, sort_keys = True)
226     labels, features = targetFeatureSplit(data)
227
228     features_train, features_test, labels_train, labels_test = train_test_split(f
229     print ""
230
231     # Criando Min/Max Scaler
232     from sklearn import preprocessing
233     scaler = preprocessing.MinMaxScaler()
234     # Scale Features
235     features = scaler.fit_transform(features)
236
237     skbest = SelectKBest(k=10) # try best value to fit
238     sk_trans = skbest.fit_transform(features_train, labels_train)
239     indices = skbest.get_support(True)
240
241     print "="*10, "skbest.scores_", "="*10
242     print skbest.scores_
243     print "="*10, "="*(len("skbest.scores_")-2), "="*10
244     print ""
245
246     print "="*10, "features - score", "="*10
247     for index in indices:
248         print 'features: %s score: %f' % (features_list[index + 1], skbest.scores_[in
249
250     print "="*10, "="*(len('features: %s score: %f')-2), "="*10

```

```

250
251 print ""
252
253 #print "GaussianNB"
254 # GaussianNB
255 clf = GaussianNB()
256 clf.fit(features_train, labels_train)
257 prediction = clf.predict(features_test)
258 print "Accuracy GaussianNB = {:.5f}".format(accuracy_score(prediction, labels_test))
259
260 #print "KNeighborsClassifier"
261 # KNeighborsClassifier
262 clf = KNeighborsClassifier()
263 clf = KNeighborsClassifier(algorithm = 'auto', leaf_size = 20, n_neighbors = 3, weight = 'distance')
264 clf.fit(features_train, labels_train)
265 prediction = clf.predict(features_test)
266 print "Accuracy KNeighborsClassifier = {:.5f}".format(accuracy_score(prediction, labels_test))
267
268 #print "SVC"
269 # SVC
270 clf = SVC(kernel = 'linear', max_iter = 10000, random_state = 42)
271 clf.fit(features_train, labels_train)
272 prediction = clf.predict(features_test)
273 print "Accuracy SVC = {:.5f}".format(accuracy_score(prediction, labels_test))
274
275 #print "AdaBoostClassifier"
276 clf = AdaBoostClassifier(DecisionTreeClassifier(max_depth=1, min_samples_leaf=2, min_samples_split=2,
277                                                    n_estimators=50, learning_rate=.8))
278 clf.fit(features_train, labels_train)
279 prediction = clf.predict(features_test)
280 print "Accuracy AdaBoostClassifier = {:.5f}".format(accuracy_score(prediction, labels_test))
281
282 def adaboost_kbest(kbest_value, new_features, features_list):
283     #Feature List
284     features_list = monta_feature(features_list)
285     if new_features == False:
286         print "Default Features"
287     # print features_list
288
289     else:
290         features_list = nova_feature(data_dict_wo, features_list)
291         print "New Features"
292     # print features_list
293     # print "Testando novos features adicionados:\n"
294     # testa_nova_feature(data_dict_wo, "DIETRICH JANET R", features_list[-2:])
295
296     # Extrair as features e os labels do conjunto de dados
297     data = featureFormat(my_dataset, features_list, sort_keys = True)
298     labels, features = targetFeatureSplit(data)
299
300     features_train, features_test, labels_train, labels_test = train_test_split(features, labels, test_size = 0.3,
301                                         random_state = 42)
302     print ""
303
304     # Criando Min/Max Scaler
305     from sklearn import preprocessing
306     scaler = preprocessing.MinMaxScaler()
307     # Scale Features
308     features = scaler.fit_transform(features)
309
310     skbest = SelectKBest(k=kbest_value) # try best value to fit
311     sk_trans = skbest.fit_transform(features_train, labels_train)
312     indices = skbest.get_support(True)

```

```

311
312
313     clf = AdaBoostClassifier(DecisionTreeClassifier(max_depth=1, min_samples_leaf=3,
314                                                     n_estimators=50, learning_rate=.95)
315
316     # Validate model precision, recall and F1-score
317     test_classifier(clf, my_dataset, features_list)
318
319
320
321 # ### Task 1: Select what features you'll use.
322 # - features_list is a list of strings, each of which is a feature name.
323 # - The first feature must be "poi".
324
325 # In[236]:
326
327
328 features_list = []
329 features_list = monta_feature(features_list)
330
331
332 # In[237]:
333
334
335 print features_list
336
337
338 # In[238]:
339
340
341 #Feature List
342 #features_list = ['poi',
343 #                 'salary',
344 #                 'deferral_payments',
345 #                 'total_payments',
346 #                 'loan_advances',
347 #                 'bonus',
348 #                 'restricted_stock_deferred',
349 #                 'deferred_income',
350 #                 'total_stock_value',
351 #                 'expenses',
352 #                 'exercised_stock_options',
353 #                 'other',
354 #                 'long_term_incentive',
355 #                 'restricted_stock',
356 #                 'director_fees',
357 #                 'to_messages',
358 #                 'from_poi_to_this_person',
359 #                 'from_messages',
360 #                 'from_this_person_to_poi',
361 #                 'shared_receipt_with_poi'] # You will need to use more features
362
363
364 # In[239]:
365
366
367 # Carregando o conjunto de dados
368 with open("final_project_dataset.pkl", "r") as data_file:
369 #     data_dict = pickle.load(data_file)
370     data_dict = pickle.load(open("final_project_dataset.pkl", "r"))
371
372

```

```
372
373
374 # In[240]:
375
376
377 print "\nQuantidade total de registros: {}\nQuantidade total de features: {}".format(
378 print "Quantidade de POI's: {}".format(conta_poi(data_dict))
379 print "Quantidade de não POI's: {}".format(len(data_dict) - conta_poi(data_dict))
380
381
382 # In[241]:
383
384
385 df_enron = pd.DataFrame.from_dict(data_dict, orient = 'index')
386
387
388 # In[242]:
389
390
391 df_enron.head()
392
393
394 # In[243]:
395
396
397 df_enron.sample(5)
398
399
400 # In[244]:
401
402
403 df_enron.tail()
404
405
406 # In[245]:
407
408
409 df_enron.describe().transpose()
410
411
412 # In[246]:
413
414
415 monta_grafico("salary", "bonus", "Primeira Analise",data_dict,1)
416
417
418 # ### Analise Exploratória
419 # - Visivelmente já é possível ver que muitas informações estão faltando
420 # - Isso fica bem evidente quando analisamos por código
421 # - Foi identificado o outliers TOTAL, que nos levou a investigar possíveis pois que
422
423 # In[247]:
424
425
426 print "Dados: HAUG DAVID L:\n\n{}".format(data_dict["TOTAL"])
427 print "Dados: LOCKHART EUGENE E:\n\n{}".format(data_dict["LOCKHART EUGENE E"])
428 print "Dados: THE TRAVEL AGENCY IN THE PARK:\n\n{}".format(data_dict["THE TRAVEL AGENI
```

AWESOME

Bom trabalho ao identificar e remover estes outliers!

Bom trabalho ao identificar e remover estes outliers!

```
429
430
431 # ### Task 2: Remove outliers
432
433 # In[248]:
434
435
436 # Removendo outliers
437 with open("final_project_dataset.pkl", "r") as data_file:
438     data_dict_woo = pickle.load(data_file)
439
440 print "Quantidade total de registros Com outliers: {}".format(len(data_dict_woo))
441
442 data_dict_woo.pop('TOTAL', None) #Não é um funcionário
443 data_dict_woo.pop('LOCKHART EUGENE E', None) #Não é um funcionário
444 data_dict_woo.pop('THE TRAVEL AGENCY IN THE PARK', None) #Não é um funcionário
445
446 print "Quantidade total de registros Sem outliers: {}".format(len(data_dict_woo))
447
448
449 # In[249]:
450
451
452 monta_grafico("salary", "bonus", "Com Outliers",data_dict,1)
453 monta_grafico("salary", "bonus", "Sem Outliers",data_dict_woo,2)
454
455
456 # In[250]:
457
458
459 df_enron_woo = pd.DataFrame.from_dict(data_dict_woo, orient = 'index')
460 df_enron_woo.describe().transpose()
461
462
463 # In[251]:
464
465
466 df_enron_woo.head(20)
467
468
469 # In[252]:
470
471
472 sumariza_valores(data_dict_woo)
473
474
475 # In[253]:
476
477
478 #df_graph_enron = df_enron_woo.copy()
479 df_graph_enron = df_enron_woo.describe().transpose()
480
481
482 # In[254]:
483
484
485 df_graph_enron.info()
486
487
```

```
488 # In[ ]:
489
490
491
492
493
494 # In[255]:
495
496
497 # Processo de criação de coluna com a diferença entre o total de registros unicos
498 #e registros duplicados/Não atribuídos
499
500 nan_dupl_col = []
501
502
503 def nan_dupl(reg):
504     qty_nan_dupl = reg['count'] - reg['unique']
505     nan_dupl_col.append(qty_nan_dupl)
506
507 df_graph_enron.apply(nan_dupl, axis=1)
508 df_graph_enron['nan_dupl'] = nan_dupl_col
509
510
511 # In[256]:
512
513
514 df_graph_enron.drop(['count', 'top', 'freq'], axis=1, inplace=True)
515
516 df_graph_enron.head()
517
518
519 # In[257]:
520
521
522 df_graph_enron.info()
523
524
525 # In[258]:
526
527
528 df_graph_enron.plot.bar()
529
530
531 # In[259]:
532
533
534 labels_graph = list(df_graph_enron.index.values)
535 print labels_graph
536
537
538 # In[260]:
539
540
541
542 bar_1 = df_graph_enron['unique']
543 bar_2 = df_graph_enron['nan_dupl']
544 x_pos = np.arange(len(bar_1))
545
546 plt2 = plt
547
548 plt.figure(figsize=(15,8))
```

```

549
550 #plt.rcParams["figure.figsize"] = [15,8]
551 #plt.rcParams["legend.frameon"] = True
552 #plt.rcParams["legend.handletextpad"] = 1
553 #plt.rcParams["legend.borderaxespad"] = 60
554
555 first_bar = plt.bar(x_pos, bar_1, 0.5, color='green')
556 second_bar = plt.bar(x_pos, bar_2, 0.5, color='skyblue', bottom=bar_1)
557 plt.title('Quantity Unique & Not a Number or Duplicated X Feature', fontsize=22)
558 plt.xlabel('Features', fontsize=18)
559 plt.ylabel('Quantity', fontsize=18)
560
561 # Definir posição e labels no eixo X
562 plt.xticks(x_pos, (labels_graph), rotation=90)
563
564 def autolabel(rects, xpos='center'):
565     """
566     Attach a text label above each bar in *rects*, displaying its height.
567
568     *xpos* indicates which side to place the text w.r.t. the center of
569     the bar. It can be one of the following {'center', 'right', 'left'}.
570     """
571
572     xpos = xpos.lower() # normalize the case of the parameter
573     ha = {'center': 'center', 'right': 'left', 'left': 'right'}
574     offset = {'center': 0.5, 'right': 0.57, 'left': 0.43} # x_txt = x + w*off
575
576
577     if rects == first_bar:
578         for rect in rects:
579             height = rect.get_height()
580             plt.text(rect.get_x() + rect.get_width()*offset[xpos], 1.00*height,
581                     '{}'.format(height), ha=ha[xpos], va='bottom')
582     else:
583         for rect in rects:
584             height = rect.get_height()
585             plt.text(rect.get_x() + rect.get_width()*offset[xpos], 1.00*138.5,
586                     '{}'.format(height), ha=ha[xpos], va='bottom')
587
588
589 autolabel(first_bar, "center")
590 autolabel(second_bar, "center")
591
592 plt.legend(labels=['Unique', 'Not a Number or Duplicated'], loc=8, borderaxespad = 47)
593
594 plt.show()
595
596
597 # In[ ]:
598
599
600
601
602
603 # ### Task 3: Create new feature(s)
604
605 # In[283]:
606
607
608 # Salvando o conjunto de dados
609 #df_enron woo = pd.DataFrame.from_dict(data_dict woo, orient = 'index')

```

```
610 #df_enron_woo.replace(to_replace='NaN', value=0.0, inplace=True)
611 my_dataset = data_dict_woo
612 #my_dataset = df_enron_woo.to_dict('index')
613
614
615 # In[284]:
616
617
618 #Feature List
619 features_list = ['poi',
620                 'salary',
621                 'deferral_payments',
622                 'total_payments',
623                 'loan_advances',
624                 'bonus',
625                 'restricted_stock_deferred',
626                 'deferred_income',
627                 'total_stock_value',
628                 'expenses',
629                 'exercised_stock_options',
630                 'other',
631                 'long_term_incentive',
632                 'restricted_stock',
633                 'director_fees',
634                 'to_messages',
635                 'from_poi_to_this_person',
636                 'from_messages',
637                 'from_this_person_to_poi',
638                 'shared_receipt_with_poi'] # You will need to use more features
639
640 features_list = nova_feature(data_dict_woo, features_list)
641 print features_list
642
643
644 # In[285]:
645
646
647 print "Testando novos features adicionados:\n"
648 testa_nova_feature(data_dict_woo, "DIETRICH JANET R", features_list[-2:])
649
650
651 # In[286]:
652
653
654 # Extraindo as features e os labels do conjunto de dados
655 data = featureFormat(my_dataset, features_list, sort_keys = True)
656 labels, features = targetFeatureSplit(data)
657
658
659 # In[287]:
660
661
662 # Criando Min/Max Scaler
663 from sklearn import preprocessing
664 scaler = preprocessing.MinMaxScaler()
665 # Scale Features
666 features = scaler.fit_transform(features)
667
668
669 # ### Task 4: Try a variety of classifiers
670 # - Please name your classifier clf for easy export below.
```

```
671 # - Note that if you want to do PCA or other multi-stage operations,  
672 # - you'll need to use Pipelines. For more info: http://scikit-learn.org/stable/modules/pipeline.html  
673  
674 # In[288]:  
675  
676  
677 #features_train, features_test, labels_train, labels_test = \  
678 #     train_test_split(features, labels, test_size=0.3, random_state=42)  
679  
680  
681 # In[289]:  
682  
683  
684 #skbest = SelectKBest(k=10) # try best value to fit  
685 #sk_trans = skbest.fit_transform(features_train, labels_train)  
686 #indices = skbest.get_support(True)  
687 #print skbest.scores_  
688  
689  
690 # In[290]:  
691  
692  
693 #for index in indices:  
694 #     print 'features: %s score: %f' % (features_list[index + 1], skbest.scores_[index])  
695 #print ""  
696  
697  
698 # In[291]:  
699  
700  
701 ##print "GaussianNB"  
702 ## GaussianNB  
703 #clf = GaussianNB()  
704 #clf.fit(features_train, labels_train)  
705 #prediction = clf.predict(features_test)  
706 #print ("Accuracy GaussianNB =", accuracy_score(prediction, labels_test))  
707  
708  
709 # In[292]:  
710  
711  
712 #print "KNeighborsClassifier"  
713 # KNeighborsClassifier  
714 #clf = KNeighborsClassifier()  
715 #clf = KNeighborsClassifier(algorithm = 'auto',leaf_size = 20,n_neighbors = 3,weights  
716 #clf.fit(features_train, labels_train)  
717 #prediction = clf.predict(features_test)  
718 #print "Accuracy KNeighborsClassifier =", accuracy_score(prediction, labels_test)  
719  
720  
721 # In[293]:  
722  
723  
724 #print "SVC"  
725 # SVC  
726 #clf = SVC(kernel = 'linear',max_iter = 10000,random_state = 42)  
727 #clf.fit(features_train, labels_train)  
728 #prediction = clf.predict(features_test)  
729 #print "Accuracy SVC =", accuracy_score(prediction, labels_test)  
730  
731
```

```

732 # In[294]:
733
734
735 #print "AdaBoostClassifier"
736 #clf = AdaBoostClassifier(DecisionTreeClassifier(max_depth=1, min_samples_leaf=2, cla:
737 #                               n_estimators=50, learning_rate=.8)
738 #clf.fit(features_train, labels_train)
739 #prediction = clf.predict(features_test)
740 #print "Accuracy AdaBoostClassifier =", accuracy_score(prediction, labels_test)
741
742
743 # In[295]:
744
745
746 accuracy(False, features_list)
747
748
749 # In[296]:
750
751
752 accuracy(True, features_list)
753
754
755 # In[297]:
756
757
758 data = {"Algorithms":["GaussianNB","KNeighborsClassifier",
759 "SVC","AdaBoostClassifier",
760 "GaussianNB","KNeighborsClassifier",
761 "SVC","AdaBoostClassifier"],
762 "New Features":["N","N","N","N","S","S","S","S"],
763 "Accuracy":[0.88372,0.90698,0.88372,0.81395,0.88372,0.90698,0.88372,0.95349],
764 }
765 #Accuracy GaussianNB = 0.88372
766 #Accuracy KNeighborsClassifier = 0.90698
767 #Accuracy SVC = 0.88372
768 #Accuracy AdaBoostClassifier = 0.81395
769 #Accuracy GaussianNB = 0.88372
770 #Accuracy KNeighborsClassifier = 0.90698
771 #Accuracy SVC = 0.88372
772 #Accuracy AdaBoostClassifier = 0.95349
773 algorithms = pd.DataFrame(data, columns = ["Algorithms", "New Features", "Accuracy",]
774 algorithms
775
776
777 # ### Task 5: Tune your classifier
778 #     Tune your classifier to achieve better than .3 precision and recall
779 #     using our testing script. Check the tester.py script in the final project
780 #     folder for details on the evaluation method, especially the test_classifier
781 #     function. Because of the small size of the dataset, the script uses
782 #     stratified shuffle split cross validation. For more info:
783
784 # In[212]:
785
786
787 # Testar os classificadores
788 #print '\n'
789 #print '##### Testar and Tuning Classifiers #####'
790 # See "pipeline_classificador" for MinMaxScaling, SelectKBest and Logistic Regression
791
792 # Classifiers tested but not using - Logistic Regression, RandomForestClassifier, Dec:

```

```

792 .....
793
794 #cross_val = pipeline_classificador('randomforest',9, features_list)
795 #print 'Melhores Parametros: ', cross_val.best_params_
796 #clf = cross_val.best_estimator_
797
798
799
800 # In[ ]:
801
802
803
804
805
806 # In[281]:
807
808
809 clf = AdaBoostClassifier(DecisionTreeClassifier(max_depth=1, min_samples_leaf=2, clas
810                                     n_estimators=50, learning_rate=.8)
811
812 #clf.fit(features_train, labels_train)
813 #prediction = clf.predict(features_test)
814 #print "Accuracy AdaBoostClassifier =", accuracy_score(prediction, labels_test)
815
816
817 # Validate model precision, recall and F1-score
818 test_classifier(clf, my_dataset, features_list)
819
820
821 # In[308]:
822
823
824 for k in range(1,11):
825     print "===== kbest = ", k, " "*10, "Ini"
826
827     t0 = time()
828     adaboost_kbest(k, True, features_list)
829
830     print "tempo de treinamento:", round(time()-t0, 3), "s"
831     print "===== kbest = ", k, " "*10, "Fim"
832
833
834
835
836 # In[309]:
837
838
839 for k in range(1,11):
840     print "===== kbest = ", k, " "*10, "Ini"
841     t0 = time()
842     adaboost_kbest(k, False, features_list)
843     print "tempo de treinamento:", round(time()-t0, 3), "s"
844     print "===== kbest = ", k, " "*10, "Fim"
845
846
847 # ### Task 6: Dump your classifier, dataset, and features_list
848 # Dump your classifier, dataset, and features_list so anyone can check your results.
849 # You do not need to change anything below, but make sure that the version of
850 # poi_id.py that you submit can be run on its own and generates the necessary .pkl
851 # files for validating your results.
852
853 # In[103]:

```

```
854
855
856 dump_classifier_and_data(clf, my_dataset, features_list)
857
858
859 # ### Referências
860 # https://datascience.stackexchange.com/questions/13410/parameters-in-gridsearchcv-in
861 #
862 # http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier
863 #
864 # https://stackoverflow.com/questions/45444953/parameter-values-for-parameter-n-estim
865 #
866 # http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.AdaBoostClassifier
867 #
868 # http://scikit-learn.org/stable/modules/generated/sklearn.model_selection.Stratified
869 #
870 # https://stackoverflow.com/questions/45969390/difference-between-stratifiedkfold-and
871 #
872 # https://www.featurelabs.com/blog/feature-engineering-vs-feature-selection/
873 #
874 # http://scikit-learn.org/stable/auto_examples/model_selection/plot_underfitting_overf
875 #
876 # http://scikit-learn.org/stable/auto_examples/feature_selection/plot_rfe_with_cross_v
877 #
878 # https://en.wikipedia.org/wiki/Enron
879 #
880 # http://scikit-learn.org/stable/modules/generated/sklearn.pipeline.Pipeline.html
881 #
882 # http://scikit-learn.org/stable/modules/generated/sklearn.grid_search.GridSearchCV.ht
883 #
884 # http://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.MinMaxScaler
885 #
886 # http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.PCA.html
887 #
888 # http://scikit-learn.org/stable/modules/generated/sklearn.feature_selection.SelectKBe
889 #
890 # http://scikit-learn.org/stable/auto_examples/plot_compare_reduction.html#sphx-glr-au
891 #
892 # https://en.wikipedia.org/wiki/Precision_and_recall
893 #
894 # http://scikit-learn.org/stable/tutorial/statistical_inference/putting_together.html
895 #
896 # https://datascience.stackexchange.com/questions/21877/how-to-use-the-output-of-grid
897 #
898 # https://stats.stackexchange.com/questions/62621/recall-and-precision-in-classificat
899 #
900 # https://stackoverflow.com/questions/49147774/what-is-random-state-in-sklearn-model-
901 #
902 # @Jefferson Aparecido Rodrigues ( Aluno Udacity que me deu umas dicas, me destravando
903
904 # In[ ]:
905
906
907
908
909
910 # In[ ]:
911
912
913
914
```


915

RETURN TO PATH
