

## File: ./pca/eigenfaces.py

```
1  """
2  =====
3  Faces recognition example using eigenfaces and SVMs
4  =====
5
6  The dataset used in this example is a preprocessed excerpt of the
7  "Labeled Faces in the Wild", aka LFW_:
8
9  http://vis-www.cs.umass.edu/lfw/lfw-funneled.tgz (233MB)
10
11  .._LFW: http://vis-www.cs.umass.edu/lfw/
12
13  original source: http://scikit-learn.org/stable/auto_examples/applications/face_recognition.html
14
15  """
16
17
18
19  print(__doc__)
20
21  from time import time
22  import logging
23  import pylab as pl
24  import numpy as np
25
26  from sklearn.model_selection import train_test_split
27  from sklearn.model_selection import GridSearchCV
28  from sklearn.datasets import fetch_lfw_people
29  from sklearn.metrics import classification_report
30  from sklearn.metrics import confusion_matrix
31  from sklearn.decomposition import PCA
32  from sklearn.svm import SVC
33
34  # Display progress logs on stdout
35  logging.basicConfig(level=logging.INFO, format='%(asctime)s %(message)s')
36
37
38  #####
39  # Download the data, if not already on disk and load it as numpy arrays
40  lfw_people = fetch_lfw_people(min_faces_per_person=70, resize=0.4)
41
42  # introspect the images arrays to find the shapes (for plotting)
43  n_samples, h, w = lfw_people.images.shape
44  np.random.seed(42)
45
46  # for machine learning we use the data directly (as relative pixel
47  # position info is ignored by this model)
48  X = lfw_people.data
49  n_features = X.shape[1]
50
51  # the label to predict is the id of the person
52  y = lfw_people.target
53  target_names = lfw_people.target_names
54  n_classes = target_names.shape[0]
55
56  print("Total dataset size:")
57  print("n_samples: %d" % n_samples)
58  print("n_features: %d" % n_features)
59  print("n_classes: %d" % n_classes)
60
61
62  #####
63  # Split into a training and testing set
64  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=42)
65
66  #####
67  # Compute a PCA (eigenfaces) on the face dataset (treated as unlabeled
68  # dataset): unsupervised feature extraction / dimensionality reduction
69  n_components = 150
70
71  print("Extracting the top %d eigenfaces from %d faces" % (n_components, X_train.shape[0]))
72  t0 = time()
73  pca = PCA(n_components=n_components, whiten=True, svd_solver='randomized').fit(X_train)
74  print("done in %0.3fs" % (time() - t0))
75
76  eigenfaces = pca.components_.reshape((n_components, h, w))
77
78  print("Projecting the input data on the eigenfaces orthonormal basis")
79  t0 = time()
```

```

80 X_train_pca = pca.transform(X_train)
81 X_test_pca = pca.transform(X_test)
82 print("done in %0.3fs" % (time() - t0))
83
84
85 #####
86 # Train a SVM classification model
87
88 print("Fitting the classifier to the training set")
89 t0 = time()
90 param_grid = {
91     'C': [1e3, 5e3, 1e4, 5e4, 1e5],
92     'gamma': [0.0001, 0.0005, 0.001, 0.005, 0.01, 0.1],
93 }
94 # for sklearn version 0.16 or prior, the class_weight parameter value is 'auto'
95 clf = GridSearchCV(SVC(kernel='rbf', class_weight='balanced'), param_grid)
96 clf = clf.fit(X_train_pca, y_train)
97 print("done in %0.3fs" % (time() - t0))
98 print("Best estimator found by grid search:")
99 print(clf.best_estimator_)
100
101
102 #####
103 # Quantitative evaluation of the model quality on the test set
104
105 print("Predicting the people names on the testing set")
106 t0 = time()
107 y_pred = clf.predict(X_test_pca)
108 print("done in %0.3fs" % (time() - t0))
109
110 print(classification_report(y_test, y_pred, target_names=target_names))
111 print(confusion_matrix(y_test, y_pred, labels=range(n_classes)))
112
113
114 #####
115 # Qualitative evaluation of the predictions using matplotlib
116
117 def plot_gallery(images, titles, h, w, n_row=3, n_col=4):
118     """Helper function to plot a gallery of portraits"""
119     plt.figure(figsize=(1.8 * n_col, 2.4 * n_row))
120     plt.subplots_adjust(bottom=0, left=.01, right=.99, top=.90, hspace=.35)
121     for i in range(n_row * n_col):
122         plt.subplot(n_row, n_col, i + 1)
123         plt.imshow(images[i].reshape((h, w)), cmap=plt.cm.gray)
124         plt.title(titles[i], size=12)
125         plt.xticks(())
126         plt.yticks(())
127
128
129 # plot the result of the prediction on a portion of the test set
130
131 def title(y_pred, y_test, target_names, i):
132     pred_name = target_names[y_pred[i]].rsplit(' ', 1)[-1]
133     true_name = target_names[y_test[i]].rsplit(' ', 1)[-1]
134     return 'predicted: %s\ntrue:      %s' % (pred_name, true_name)
135
136 prediction_titles = [title(y_pred, y_test, target_names, i)
137                      for i in range(y_pred.shape[0])]
138
139 plot_gallery(X_test, prediction_titles, h, w)
140
141
142 # plot the gallery of the most significant eigenfaces
143 eigenface_titles = ["eigenface %d" % i for i in range(eigenfaces.shape[0])]
144 plot_gallery(eigenfaces, eigenface_titles, h, w)
145
146 plt.show()

```

## File: ./naive\_bayes/nb\_author\_id.py

```

1 #!/usr/bin/python
2
3 """
4 This is the code to accompany the Lesson 1 (Naive Bayes) mini-project.
5
6 Use a Naive Bayes Classifier to identify emails by their authors
7
8 authors and labels:

```

```

9   Sara has label 0
10  Chris has label 1
11  """
12
13  import sys
14  from time import time
15  sys.path.append("./tools/")
16  from email_preprocess import preprocess
17
18
19  #### features_train and features_test are the features for the training
20  #### and testing datasets, respectively
21  #### labels_train and labels_test are the corresponding item labels
22  features_train, features_test, labels_train, labels_test = preprocess()
23
24
25
26
27  #####
28  ### your code goes here ###
29
30
31  #####
32
33

```

### File: ./tools/feature\_format.py

```

1  #!/usr/bin/python
2
3  """
4      A general tool for converting data from the
5      dictionary format to an (n x k) python list that's
6      ready for training an sklearn algorithm
7
8      n--no. of key-value pairs in dictionary
9      k--no. of features being extracted
10
11     dictionary keys are names of persons in dataset
12     dictionary values are dictionaries, where each
13     key-value pair in the dict is the name
14     of a feature, and its value for that person
15
16     In addition to converting a dictionary to a numpy
17     array, you may want to separate the labels from the
18     features--this is what targetFeatureSplit is for
19
20     so, if you want to have the poi label as the target,
21     and the features you want to use are the person's
22     salary and bonus, here's what you would do:
23
24     feature_list = ["poi", "salary", "bonus"]
25     data_array = featureFormat( data_dictionary, feature_list )
26     label, features = targetFeatureSplit(data_array)
27
28     the line above (targetFeatureSplit) assumes that the
29     label is the _first_ item in feature_list--very important
30     that poi is listed first!
31  """
32
33
34  import numpy as np
35
36  def featureFormat( dictionary, features, remove_NaN=True, remove_all_zeroes=True, remove_any_zeroes=False, sort_keys = False):
37      """ convert dictionary to numpy array of features
38      remove_NaN = True will convert "NaN" string to 0.0
39      remove_all_zeroes = True will omit any data points for which
40      all the features you seek are 0.0
41      remove_any_zeroes = True will omit any data points for which
42      any of the features you seek are 0.0
43      sort_keys = True sorts keys by alphabetical order. Setting the value as
44      a string opens the corresponding pickle file with a preset key
45      order (this is used for Python 3 compatibility, and sort_keys
46      should be left as False for the course mini-projects).
47      NOTE: first feature is assumed to be 'poi' and is not checked for
48      removal for zero or missing values.
49  """
50
51
52  return list = []

```

```

52     return_list = []
53
54     # Key order - first branch is for Python 3 compatibility on mini-projects,
55     # second branch is for compatibility on final project.
56     if isinstance(sort_keys, str):
57         import pickle
58         keys = pickle.load(open(sort_keys, "rb"))
59     elif sort_keys:
60         keys = sorted(dictionary.keys())
61     else:
62         keys = dictionary.keys()
63
64     for key in keys:
65         tmp_list = []
66         for feature in features:
67             try:
68                 dictionary[key][feature]
69             except KeyError:
70                 print("error: key ", feature, " not present")
71                 return
72             value = dictionary[key][feature]
73             if value=="NaN" and remove_NaN:
74                 value = 0
75             tmp_list.append( float(value) )
76
77         # Logic for deciding whether or not to add the data point.
78         append = True
79         # exclude 'poi' class as criteria.
80         if features[0] == 'poi':
81             test_list = tmp_list[1:]
82         else:
83             test_list = tmp_list
84         ### if all features are zero and you want to remove
85         ### data points that are all zero, do that here
86         if remove_all_zeroes:
87             append = False
88             for item in test_list:
89                 if item != 0 and item != "NaN":
90                     append = True
91                     break
92         ### if any features for a given data point are zero
93         ### and you want to remove data points with any zeroes,
94
95         ### handle that here
96         if remove_any_zeroes:
97             if 0 in test_list or "NaN" in test_list:
98                 append = False
99         ### Append the data point if flagged for addition.
100        if append:
101            return_list.append( np.array(tmp_list) )
102
103    return np.array(return_list)
104
105 def targetFeatureSplit( data ):
106     """
107     given a numpy array like the one returned from
108     featureFormat, separate out the first feature
109     and put it into its own list (this should be the
110     quantity you want to predict)
111
112     return targets and features as separate lists
113
114     (sklearn can generally handle both lists and numpy arrays as
115     input formats when training/predicting)
116     """
117
118     target = []
119     features = []
120     for item in data:
121         target.append( item[0] )
122         features.append( item[1:] )
123
124     return target, features
125
126
127
128

```

**File: ./tools/startup.py**

```

1  #!/usr/bin/python
2  import time
3  import sys
4
5  if sys.version_info[0] >= 3:
6      from urllib.request import urlretrieve
7  else:
8      # Not Python 3 - today, it is most likely to be Python 2
9      # But note that this might need an update when Python 4
10     # might be around one day
11     from urllib import urlretrieve
12
13 print()
14 print("checking for nltk")
15 try:
16     import nltk
17 except ImportError:
18     print("you should install nltk before continuing")
19
20 print("checking for numpy")
21 try:
22     import numpy
23 except ImportError:
24     print("you should install numpy before continuing")
25
26 print("checking for scipy")
27 try:
28     import scipy
29 except:
30     print("you should install scipy before continuing")
31
32 print("checking for sklearn")
33 try:
34     import sklearn
35 except:
36     print("you should install sklearn before continuing")
37
38 print()
39 print("downloading the Enron dataset (this may take a while)")
40 print("to check on progress, you can cd up one level, then execute <ls -lthr>")
41 print("Enron dataset should be last item on the list, along with its current size")
42 print("download will complete at about 423 MB")
43
44
45 def reporthook(count, block_size, total_size):
46     """
47     Callback for displaying progress bar while downloading the enron dataset.
48     Thanks to Pushkaraj Shinde https://github.com/udacity/ud120-projects/pull/55
49     For further information about parameters see documentation of urlretrieve,
50     e.g. https://docs.python.org/3.5/library/urllib.request.html#urllib.request.urlretrieve
51     """
52     global start_time
53     if count == 0:
54         start_time = time.time()
55     return
56     duration = time.time() - start_time
57     progress_size = int(count * block_size)
58     speed = int(progress_size / (1024 * duration))
59     percent = min(int(count * block_size * 100 / total_size), 100)
60     sys.stdout.write("\rEnron dataset: Downloaded %d%%, %d MB, %d KB/s, %d seconds passed" %
61                     (percent, progress_size / (1024 * 1024), speed, duration))
62     sys.stdout.flush()
63
64 url = "https://www.cs.cmu.edu/~./enron/enron_mail_20150507.tgz"
65
66 urlretrieve(url, "../enron_mail_20150507.tgz", reporthook)
67 print("download complete!")
68
69
70 print()
71 print("unzipping Enron dataset (this may take a while)")
72 import tarfile
73 import os
74 os.chdir("..")
75 tfile = tarfile.open("enron_mail_20150507.tar.gz", "r:gz")
76 tfile.extractall(".")
77
78 print("you're ready to go!")

```

**File: ./tools/parse\_out\_email\_text.py**

```
1 #!/usr/bin/python
2
3 from nltk.stem.snowball import SnowballStemmer
4 import string
5 import re
6
7 def parseOutText(f):
8     """ given an opened email file f, parse out all text below the
9     metadata block at the top
10     (in Part 2, you will also add stemming capabilities)
11     and return a string that contains all the words
12     in the email (space-separated)
13
14     example use case:
15     f = open("email_file_name.txt", "r")
16     text = parseOutText(f)
17
18     """
19
20
21 f.seek(0) ### go back to beginning of file (annoying)
22 all_text = f.read()
23
24 ### split off metadata
25 content = all_text.split("X-FileName:")
26 words = ""
27 if len(content) > 1:
28     ### remove punctuation
29     text_string = re.sub('[! + string.punctuation + \']', "", content[1])
30     text_string = content[1].translate(str.maketrans("", "", string.punctuation))
31
32     ### project part 2: comment out the line below
33     #words = text_string
34
35     ### split the text string into individual words, stem each word,
36     ### and append the stemmed word to words (make sure there's a single
37     ### space between each stemmed word)
38     stemmer = SnowballStemmer("english")
39     word_list = []
40
41     for word in text_string.split():
42         word_list.append(stemmer.stem(word))
43
44     words = ' '.join(word_list)
45 return words
46
47
48
49 def main():
50     ff = open("../text_learning/test_email.txt", "r")
51     text = parseOutText(ff)
52     print(text)
53
54
55
56 if __name__ == '__main__':
57     main()
58
```

**File: ./tools/email\_preprocess.py**

```
1 #!/usr/bin/python
2
3 import pickle
4 #import cPickle
5 import numpy
6
7 from sklearn.feature_extraction.text import TfidfVectorizer
8 from sklearn.feature_selection import SelectPercentile, f_classif
9 from sklearn.model_selection import train_test_split
10
11
12 def preprocess(words_file = "../tools/word_data.pkl", authors_file = "../tools/email_authors.pkl"):
13     """
14     this function takes a pre-made list of email texts (by default word_data.pkl)
15     and the corresponding authors (by default email_authors.pkl) and performs

```

```

16     a number of preprocessing steps:
17     -- splits into training/testing sets (10% testing)
18     -- vectorizes into tfidf matrix
19     -- selects/keeps most helpful features
20
21     after this, the features and labels are put into numpy arrays, which play nice with sklearn functions
22
23     4 objects are returned:
24     -- training/testing features
25     -- training/testing labels
26
27     """
28
29     ### the words (features) and authors (labels), already largely preprocessed
30     ### this preprocessing will be repeated in the text learning mini-project
31     authors_file_handler = open(authors_file, "rb")
32     authors = pickle.load(authors_file_handler)
33     authors_file_handler.close()
34
35     words_file_handler = open(words_file, "rb")
36     word_data = pickle.load(words_file_handler)
37     words_file_handler.close()
38
39     ### test_size is the percentage of events assigned to the test set
40     ### (remainder go into training)
41     features_train, features_test, labels_train, labels_test = \
42         train_test_split(word_data, authors, test_size=0.1, random_state=42)
43
44
45
46     ### text vectorization--go from strings to lists of numbers
47     vectorizer = TfidfVectorizer(sublinear_tf=True, max_df=0.5,
48                                 stop_words='english')
49     features_train_transformed = vectorizer.fit_transform(features_train)
50     features_test_transformed = vectorizer.transform(features_test)
51
52
53
54     ### feature selection, because text is super high dimensional and
55     ### can be really computationally chewy as a result
56     selector = SelectPercentile(f_classif, percentile=10)
57     selector.fit(features_train_transformed, labels_train)
58     features_train_transformed = selector.transform(features_train_transformed).toarray()
59     features_test_transformed = selector.transform(features_test_transformed).toarray()
60
61     ### info on the data
62     print("no. of Chris training emails:", sum(labels_train))
63     print("no. of Sara training emails:", len(labels_train)-sum(labels_train))
64
65     return features_train_transformed, features_test_transformed, labels_train, labels_test

```

## File: ./text\_learning/vectorize\_text.py

```

1  #!/usr/bin/python
2
3  import os
4  import pickle
5  import re
6  import sys
7
8  sys.path.append( "../tools/" )
9  from parse_out_email_text import parseOutText
10
11  """
12  Starter code to process the emails from Sara and Chris to extract
13  the features and get the documents ready for classification.
14
15  The list of all the emails from Sara are in the from_sara list
16  likewise for emails from Chris (from_chris)
17
18  The actual documents are in the Enron email dataset, which
19  you downloaded/unpacked in Part 0 of the first mini-project. If you have
20  not obtained the Enron email corpus, run startup.py in the tools folder.
21
22  The data is stored in lists and packed away in pickle files at the end.
23  """
24
25
26  from_sara = open("from_sara.txt", "r")

```

```

27 from_chris = open("from_chris.txt", "r")
28
29 from_data = []
30 word_data = []
31
32 ### temp_counter is a way to speed up the development--there are
33 ### thousands of emails from Sara and Chris, so running over all of them
34 ### can take a long time
35 ### temp_counter helps you only look at the first 200 emails in the list so you
36 ### can iterate your modifications quicker
37 temp_counter = 0
38
39
40 for name, from_person in [("sara", from_sara), ("chris", from_chris)]:
41     for path in from_person:
42         ### only look at first 200 emails when developing
43         ### once everything is working, remove this line to run over full dataset
44         temp_counter += 1
45         if temp_counter < 200:
46             path = os.path.join('.', path[:-1])
47             print(path)
48             email = open(path, "r")
49
50             ### use parseOutText to extract the text from the opened email
51
52             ### use str.replace() to remove any instances of the words
53             ### ["sara", "shackleton", "chris", "germani"]
54
55             ### append the text to word_data
56
57             ### append a 0 to from_data if email is from Sara, and 1 if email is from Chris
58
59
60             email.close()
61
62 print("emails processed")
63 from_sara.close()
64 from_chris.close()
65
66 pickle.dump(word_data, open("your_word_data.pkl", "wb"))
67 pickle.dump(from_data, open("your_email_authors.pkl", "wb"))
68
69
70
71
72
73 ### in Part 4, do Tfidf vectorization here
74
75

```

#### File: ./outliers/outlier\_cleaner.py

```

1 #!/usr/bin/python
2
3
4 def outlierCleaner(predictions, ages, net_worths):
5     """
6     Clean away the 10% of points that have the largest
7     residual errors (difference between the prediction
8     and the actual net worth).
9
10    Return a list of tuples named cleaned_data where
11    each tuple is of the form (age, net_worth, error).
12    """
13
14    cleaned_data = []
15
16    ### your code goes here
17
18
19    return cleaned_data
20

```

#### File: ./outliers/enron\_outliers.py

```

1 #!/usr/bin/python

```



```

2
3 import pickle
4 import sys
5 import matplotlib.pyplot
6 sys.path.append("../tools/")
7 from feature_format import featureFormat, targetFeatureSplit
8
9
10 ### read in data dictionary, convert to numpy array
11 data_dict = pickle.load(open("../final_project/final_project_dataset.pkl", "rb"))
12 features = ["salary", "bonus"]
13 data = featureFormat(data_dict, features)
14
15
16 ### your code below
17
18
19

```

# File: ./outliers/outlier\_removal\_regression.py

```

1 #!/usr/bin/python
2
3 import random
4 import numpy
5 import matplotlib.pyplot as plt
6 import pickle
7
8 from outlier_cleaner import outlierCleaner
9
10
11 ### load up some practice data with outliers in it
12 ages = pickle.load(open("practice_outliers_ages.pkl", "rb"))
13 net_worths = pickle.load(open("practice_outliers_net_worths.pkl", "rb"))
14
15
16
17 ### ages and net_worths need to be reshaped into 2D numpy arrays
18 ### second argument of reshape command is a tuple of integers: (n_rows, n_columns)
19 ### by convention, n_rows is the number of data points
20 ### and n_columns is the number of features
21 ages = numpy.reshape(numpy.array(ages), (len(ages), 1))
22 net_worths = numpy.reshape(numpy.array(net_worths), (len(net_worths), 1))
23 from sklearn.model_selection import train_test_split
24 ages_train, ages_test, net_worths_train, net_worths_test = train_test_split(ages, net_worths, test_size=0.1, random_state=42)
25
26 ### fill in a regression here! Name the regression object reg so that
27 ### the plotting code below works, and you can see what your regression looks like
28 from sklearn.linear_model import LinearRegression
29 reg = LinearRegression()
30 reg.fit(ages_train, net_worths_train)
31 print("coefficient of determination R^2 on train data {:.3}"
32       format(reg.score(ages_test, net_worths_test)))
33
34
35 try:
36     plt.plot(ages, reg.predict(ages), color="blue")
37 except NameError:
38     pass
39 plt.scatter(ages, net_worths)
40 plt.show()
41
42
43 ### identify and remove the most outlier-y points
44 cleaned_data = []
45 try:
46     predictions = reg.predict(ages_train)
47     cleaned_data = outlierCleaner(predictions, ages_train, net_worths_train)
48 except NameError:
49     print("your regression object doesn't exist, or isn't name reg")
50     print("can't make predictions to use in identifying outliers")
51
52
53
54
55
56
57
58 ### only run this code if cleaned_data is returning data

```

```

59 if len(cleaned_data) > 0:
60     ages, net_worths, errors = zip(*cleaned_data)
61     ages = numpy.reshape(numpy.array(ages), (len(ages), 1))
62     net_worths = numpy.reshape(numpy.array(net_worths), (len(net_worths), 1))
63
64     ### refit your cleaned data!
65     try:
66         reg_fit(ages, net_worths)
67         plt.plot(ages, reg.predict(ages), color="blue")
68         print("coefficient of determination R^2 on train data {:.3}"
69               format(reg.score(ages_test, net_worths_test)))
70     except NameError:
71         print("you don't seem to have regression imported/created,")
72         print(" or else your regression object isn't named reg")
73         print(" either way, only draw the scatter plot of the cleaned data")
74     plt.scatter(ages, net_worths)
75     plt.xlabel("ages")
76     plt.ylabel("net worths")
77     plt.show()
78
79
80 else:
81     print("outlierCleaner() is returning an empty list, no refitting to be done")
82

```

### File: ./validation/validate\_poi.py

```

1  #!/usr/bin/python
2
3
4  """
5      Starter code for the validation mini-project.
6      The first step toward building your POI identifier!
7
8      Start by loading/formatting the data
9
10     After that, it's not our code anymore--it's yours!
11 """
12
13 import pickle
14 import sys
15 from sklearn.tree import DecisionTreeClassifier
16 from sklearn.model_selection import train_test_split
17 sys.path.append("../tools/")
18 from feature_format import featureFormat, targetFeatureSplit
19 from mylib import fit_and_predict
20
21 data_dict = pickle.load(open("../final_project/final_project_dataset.pkl", "rb"))
22
23 ### first element is our labels, any added elements are predictor
24 ### features. Keep this the same for the mini-project, but you'll
25 ### have a different feature list when you do the final project.
26 sort_keys = "../tools/python2_lesson13_keys.pkl"
27 features_list = ["poi", "salary"]
28
29 data = featureFormat(data_dict, features_list)
30 labels, features = targetFeatureSplit(data)
31
32 clf = DecisionTreeClassifier(random_state=42)
33 acc = fit_and_predict(clf, features, features, labels, labels)
34 print("Overfitted {:.3}" format(acc))
35
36 clf = DecisionTreeClassifier(random_state=42)
37 X_train, X_test, y_train, y_test = train_test_split(features, labels, test_size=0.30, random_state=42, shuffle=False)
38 acc = fit_and_predict(clf, X_train, X_test, y_train, y_test)
39 print("Accuracy of test {:.3}" format(acc))
40
41 ### it's all yours from here forward!
42
43

```

### File: ./code2pdf

:directories:  
- .git

:files:  
- .txt  
- .ignore  
- README.md

**File: ./svm/svm\_author\_id.py**

```
1 #!/usr/bin/python
2
3 """
4     This is the code to accompany the Lesson 2 (SVM) mini-project.
5
6     Use a SVM to identify emails from the Enron corpus by their authors:
7     Sara has label 0
8     Chris has label 1
9 """
10
11 import sys
12 from time import time
13 sys.path.append("../tools/")
14 from email_preprocess import preprocess
15
16
17 ### features_train and features_test are the features for the training
18 ### and testing datasets, respectively
19 ### labels_train and labels_test are the corresponding item labels
20 features_train, features_test, labels_train, labels_test = preprocess()
21
22
23
24
25 #####
26 ### your code goes here ###
27
28 #####
29
30
```

**File: ./small\_task/count\_stop\_words.py**

```
1 from nltk.corpus import stopwords
2
3 sw = stopwords.words("english")
4 print(len(sw))
5 print(sw)
```

**File: ./datasets\_questions/explore\_enron\_data.py**

```
1 #!/usr/bin/python
2
3 """
4     Starter code for exploring the Enron dataset (emails + finances);
5     loads up the dataset (pickled dict of dicts).
6
7     The dataset has the form:
8     enron_data["LASTNAME FIRSTNAME MIDDLEINITIAL"] = { features_dict }
9
10     {features_dict} is a dictionary of features associated with that person.
11     You should explore features_dict as part of the mini-project,
12     but here's an example to get you started:
13
14     enron_data["SKILLING JEFFREY K"]["bonus"] = 5600000
15
16 """
17
18 import pickle
19
20 enron_data = pickle.load(open("../final_project/final_project_dataset.pkl", "rb"))
21
22
```

**File: ./final\_project/poi\_id.py**

```
1  #!/usr/bin/python
2
3  import sys
4  import pickle
5  sys.path.append("../tools/")
6
7  from feature_format import featureFormat, targetFeatureSplit
8  from tester import dump_classifier_and_data
9
10 ##### Task 1: Select what features you'll use.
11 ##### features_list is a list of strings, each of which is a feature name.
12 ##### The first feature must be "poi".
13 features_list = ['poi','salary'] # You will need to use more features
14
15 ##### Load the dictionary containing the dataset
16 with open("final_project_dataset.pkl", "r") as data_file:
17     data_dict = pickle.load(data_file)
18
19 ##### Task 2: Remove outliers
20 ##### Task 3: Create new feature(s)
21 ##### Store to my_dataset for easy export below.
22 my_dataset = data_dict
23
24 ##### Extract features and labels from dataset for local testing
25 data = featureFormat(my_dataset, features_list, sort_keys = True)
26 labels, features = targetFeatureSplit(data)
27
28 ##### Task 4: Try a variety of classifiers
29 ##### Please name your classifier clf for easy export below.
30 ##### Note that if you want to do PCA or other multi-stage operations,
31 ##### you'll need to use Pipelines. For more info:
32 ##### http://scikit-learn.org/stable/modules/pipeline.html
33
34 # Provided to give you a starting point. Try a variety of classifiers.
35 from sklearn.naive_bayes import GaussianNB
36 clf = GaussianNB()
37
38 ##### Task 5: Tune your classifier to achieve better than .3 precision and recall
39 ##### using our testing script. Check the tester.py script in the final project
40 ##### folder for details on the evaluation method, especially the test_classifier
41 ##### function. Because of the small size of the dataset, the script uses
42 ##### stratified shuffle split cross validation. For more info:
43 ##### http://scikit-learn.org/stable/modules/generated/sklearn.cross_validation.StratifiedShuffleSplit.html
44
45 # Example starting point. Try investigating other evaluation techniques!
46 from sklearn.cross_validation import train_test_split
47 features_train, features_test, labels_train, labels_test = \
48     train_test_split(features, labels, test_size=0.3, random_state=42)
49
50 ##### Task 6: Dump your classifier, dataset, and features_list so anyone can
51 ##### check your results. You do not need to change anything below, but make sure
52 ##### that the version of poi_id.py that you submit can be run on its own and
53 ##### generates the necessary .pkl files for validating your results.
54
55 dump_classifier_and_data(clf, my_dataset, features_list)
```

**File: ./final\_project/tester.py**

```
1  #!/usr/bin/pickle
2
3  """ a basic script for importing student's POI identifier,
4      and checking the results that they get from it
5
6      requires that the algorithm, dataset, and features list
7      be written to my_classifier.pkl, my_dataset.pkl, and
8      my_feature_list.pkl, respectively
9
10     that process should happen at the end of poi_id.py
11 """
12
13 import pickle
14 import sys
15 from sklearn.cross_validation import StratifiedShuffleSplit
16 sys.path.append("../tools/")
17 from feature_format import featureFormat, targetFeatureSplit
18
```

```

19 PERF_FORMAT_STRING = "\
20 \tAccuracy: {:.>0.{display_precision}f}\tPrecision: {:.>0.{display_precision}f}\t\
21 Recall: {:.>0.{display_precision}f}\tF1: {:.>0.{display_precision}f}\tF2: {:.>0.{display_precision}f}"
22 RESULTS_FORMAT_STRING = "\tTotal predictions: {:.4d}\tTrue positives: {:.4d}\tFalse positives: {:.4d}\
23 \tFalse negatives: {:.4d}\tTrue negatives: {:.4d}"
24
25 def test_classifier(clf, dataset, feature_list, folds = 1000):
26     data = featureFormat(dataset, feature_list, sort_keys = True)
27     labels, features = targetFeatureSplit(data)
28     cv = StratifiedShuffleSplit(labels, folds, random_state = 42)
29     true_negatives = 0
30     false_negatives = 0
31     true_positives = 0
32     false_positives = 0
33     for train_idx, test_idx in cv:
34         features_train = []
35         features_test = []
36         labels_train = []
37         labels_test = []
38         for ii in train_idx:
39             features_train.append( features[ii] )
40             labels_train.append( labels[ii] )
41         for jj in test_idx:
42             features_test.append( features[jj] )
43             labels_test.append( labels[jj] )
44
45     ### fit the classifier using training set, and test on test set
46     clf.fit(features_train, labels_train)
47     predictions = clf.predict(features_test)
48     for prediction, truth in zip(predictions, labels_test):
49         if prediction == 0 and truth == 0:
50             true_negatives += 1
51         elif prediction == 0 and truth == 1:
52             false_negatives += 1
53         elif prediction == 1 and truth == 0:
54             false_positives += 1
55         elif prediction == 1 and truth == 1:
56             true_positives += 1
57         else:
58             print "Warning: Found a predicted label not == 0 or 1."
59             print "All predictions should take value 0 or 1."
60             print "Evaluating performance for processed predictions:"
61             break
62     try:
63         total_predictions = true_negatives + false_negatives + false_positives + true_positives
64         accuracy = 1.0*(true_positives + true_negatives)/total_predictions
65         precision = 1.0*true_positives/(true_positives+false_positives)
66         recall = 1.0*true_positives/(true_positives+false_negatives)
67         f1 = 2.0 * true_positives/(2* true_positives + false_positives+false_negatives)
68         f2 = (1+2.0*2.0) * precision*recall/(4*precision + recall)
69         print clf
70         print PERF_FORMAT_STRING.format(accuracy, precision, recall, f1, f2, display_precision = 5)
71         print RESULTS_FORMAT_STRING.format(total_predictions, true_positives, false_positives, false_negatives, true_negatives)
72         print ""
73     except:
74         print "Got a divide by zero when trying out:", clf
75         print "Precision or recall may be undefined due to a lack of true positive predicitions."
76
77 CLF_PICKLE_FILENAME = "my_classifier.pkl"
78 DATASET_PICKLE_FILENAME = "my_dataset.pkl"
79 FEATURE_LIST_FILENAME = "my_feature_list.pkl"
80
81 def dump_classifier_and_data(clf, dataset, feature_list):
82     with open(CLF_PICKLE_FILENAME, "w") as clf_outfile:
83         pickle.dump(clf, clf_outfile)
84     with open(DATASET_PICKLE_FILENAME, "w") as dataset_outfile:
85         pickle.dump(dataset, dataset_outfile)
86     with open(FEATURE_LIST_FILENAME, "w") as featurelist_outfile:
87         pickle.dump(feature_list, featurelist_outfile)
88
89 def load_classifier_and_data():
90     with open(CLF_PICKLE_FILENAME, "r") as clf_infile:
91         clf = pickle.load(clf_infile)
92     with open(DATASET_PICKLE_FILENAME, "r") as dataset_infile:
93         dataset = pickle.load(dataset_infile)
94     with open(FEATURE_LIST_FILENAME, "r") as featurelist_infile:
95         feature_list = pickle.load(featurelist_infile)
96     return clf, dataset, feature_list
97
98 def main():
99     ### load up student's classifier, dataset, and feature_list

```

```

100 clf, dataset, feature_list = load_classifier_and_data()
101 ### Run testing script
102 test_classifier(clf, dataset, feature_list)
103
104 if __name__ == '__main__':
105     main()

```

**File: ./final\_project/poi\_email\_addresses.py**

```

1 def poiEmails():
2     email_list = ["kenneth_lay@enron.net",
3                  "kenneth_lay@enron.com",
4                  "klay.enron@enron.com",
5                  "kenneth.lay@enron.com",
6                  "klay@enron.com",
7                  "layk@enron.com",
8                  "chairman.ken@enron.com",
9                  "jeffreyskilling@yahoo.com",
10                 "jeff_skilling@enron.com",
11                 "jskilling@enron.com",
12                 "effrey.skilling@enron.com",
13                 "skilling@enron.com",
14                 "jeffrey.k.skilling@enron.com",
15                 "jeff.skilling@enron.com",
16                 "kevin_a_howard.enronxgate.enron@enron.net",
17                 "kevin.howard@enron.com",
18                 "kevin.howard@enron.net",
19                 "kevin.howard@gcm.com",
20                 "michael.krautz@enron.com",
21                 "scott.yeager@enron.com",
22                 "syeager@fyi-net.com",
23                 "scott_yeager@enron.net",
24                 "syeager@flash.net",
25                 "joe'.hirko@enron.com",
26                 "joe.hirko@enron.com",
27                 "rex.shelby@enron.com",
28                 "rex.shelby@enron.nt",
29                 "rex_shelby@enron.net",
30                 "jbrown@enron.com",
31                 "james.brown@enron.com",
32                 "rick.causey@enron.com",
33                 "richard.causey@enron.com",
34                 "rcausey@enron.com",
35                 "calger@enron.com",
36                 "chris.calger@enron.com",
37                 "christopher.calger@enron.com",
38                 "ccalger@enron.com",
39                 "tim_despain.enronxgate.enron@enron.net",
40                 "tim.despain@enron.com",
41                 "kevin_hannon@enron.com",
42                 "kevin'.hannon@enron.com",
43                 "kevin_hannon@enron.net",
44                 "kevin.hannon@enron.com",
45                 "mkoenig@enron.com",
46                 "mark.koenig@enron.com",
47                 "m..forney@enron.com",
48                 "ken'.rice@enron.com",
49                 "ken.rice@enron.com",
50                 "ken_rice@enron.com",
51                 "ken_rice@enron.net",
52                 "paula.rieker@enron.com",
53                 "priecker@enron.com",
54                 "andrew.fastow@enron.com",
55                 "lfastow@pdq.net",
56                 "andrew.s.fastow@enron.com",
57                 "lfastow@pop.pdq.net",
58                 "andy.fastow@enron.com",
59                 "david.w.delainey@enron.com",
60                 "delainey.dave@enron.com",
61                 "'delainey@enron.com",
62                 "david.delainey@enron.com",
63                 "david.delainey'@enron.com",
64                 "dave.delainey@enron.com",
65                 "delainey'.david@enron.com",
66                 "ben.glisan@enron.com",
67                 "bglisan@enron.com",
68                 "ben_f_glisan@enron.com",
69                 "ben'.glisan@enron.com",
70                 "jeff.richter@enron.com",

```

```

71     "jrichter@nwlinc.com",
72     "lawrencelawyer@aol.com",
73     "lawyer'.larry@enron.com",
74     "larry_lawyer@enron.com",
75     "llawyer@enron.com",
76     "larry.lawyer@enron.com",
77     "lawrence.lawyer@enron.com",
78     "tbelden@enron.com",
79     "tim.belden@enron.com",
80     "tim_belden@pgn.com",
81     "tbelden@ect.enron.com",
82     "michael.kopper@enron.com",
83     "dave.duncan@enron.com",
84     "dave.duncan@cipco.org",
85     "duncan.dave@enron.com",
86     "ray.bowen@enron.com",
87     "raymond.bowen@enron.com",
88     "bowen@enron.com",
89     "wes.colwell@enron.com",
90     "dan.boyle@enron.com",
91     "cloehr@enron.com",
92     "chris.loehr@enron.com"
93 ]
94 return email_list

```

#### File: ./decision\_tree/dt\_author\_id.py

```

1  #!/usr/bin/python
2
3  """
4      This is the code to accompany the Lesson 3 (decision tree) mini-project.
5
6      Use a Decision Tree to identify emails from the Enron corpus by author:
7      Sara has label 0
8      Chris has label 1
9  """
10
11 import sys
12 from time import time
13 sys.path.append("../tools/")
14 from email_preprocess import preprocess
15
16
17 #### features_train and features_test are the features for the training
18 #### and testing datasets, respectively
19 #### labels_train and labels_test are the corresponding item labels
20 features_train, features_test, labels_train, labels_test = preprocess()
21
22
23
24
25 #####
26 ### your code goes here ###
27
28
29 #####
30
31

```

#### File: ./feature\_selection/find\_signature.py

```

1  #!/usr/bin/python
2
3  import pickle
4  import numpy
5  numpy.random.seed(42)
6
7
8  #### The words (features) and authors (labels), already largely processed.
9  #### These files should have been created from the previous (Lesson 10)
10 #### mini-project.
11 words_file = "../text_learning/your_word_data.pkl"
12 authors_file = "../text_learning/your_email_authors.pkl"
13 word_data = pickle.load( open(words_file, "rb"))
14 authors = pickle.load( open(authors_file, "rb") )
15

```

```

16
17
18 ### test_size is the percentage of events assigned to the test set (the
19 ### remainder go into training)
20 ### feature matrices changed to dense representations for compatibility with
21 ### classifier functions in versions 0.15.2 and earlier
22 from sklearn.model_selection import train_test_split
23 features_train, features_test, labels_train, labels_test = \
24     train_test_split(word_data, authors, test_size=0.1, random_state=42)
25
26 from sklearn.feature_extraction.text import TfidfVectorizer
27 vectorizer = TfidfVectorizer(sublinear_tf=True, max_df=0.5,
28                             stop_words='english')
29 features_train = vectorizer.fit_transform(features_train)
30 features_test = vectorizer.transform(features_test).toarray()
31
32
33 ### a classic way to overfit is to use a small number
34 ### of data points and a large number of features;
35 ### train on only 150 events to put ourselves in this regime
36 features_train = features_train[:150].toarray()
37 labels_train = labels_train[:150]
38
39
40
41 ### your code goes here
42
43
44

```

#### File: ./regression/finance\_regression.py

```

1  #!/usr/bin/python
2
3  """
4      Starter code for the regression mini-project.
5
6      Loads up/formats a modified version of the dataset
7      (why modified? we've removed some trouble points
8      that you'll find yourself in the outliers mini-project).
9
10     Draws a little scatterplot of the training/testing data
11
12     You fill in the regression code where indicated:
13 """
14
15
16 import sys
17 import pickle
18 sys.path.append("../tools/")
19 from feature_format import featureFormat, targetFeatureSplit
20 dictionary = pickle.load( open("../final_project/final_project_dataset_modified.pkl", "rb") )
21
22 ### list the features you want to look at--first item in the
23 ### list will be the "target" feature
24 features_list = ["bonus", "salary"]
25 data = featureFormat(dictionary, features_list, remove_any_zeroes=True,
26                      sort_keys='../tools/python2_lesson06_keys.pkl')
27 # see https://classroom.udacity.com/courses/ud120/lessons/2301748537/concepts/30416086000923
28 target, features = targetFeatureSplit( data )
29
30 ### training-testing split needed in regression, just like classification
31 from sklearn.model_selection import train_test_split
32 feature_train, feature_test, target_train, target_test = train_test_split(features, target, test_size=0.5, random_state=42)
33 train_color = "b"
34 test_color = "b"
35
36
37
38 ### Your regression goes here!
39 ### Please name it reg, so that the plotting code below picks it up and
40 ### plots it correctly. Don't forget to change the test_color above from "b" to
41 ### "r" to differentiate training points from test points.
42
43
44
45
46
47

```



```

48
49
50 ### draw the scatterplot, with color-coded training and testing points
51 import matplotlib.pyplot as plt
52 for feature, target in zip(feature_test, target_test):
53     plt.scatter( feature, target, color=test_color )
54 for feature, target in zip(feature_train, target_train):
55     plt.scatter( feature, target, color=train_color )
56
57 ### labels for the legend
58 plt.scatter(feature_test[0], target_test[0], color=test_color, label="test")
59 plt.scatter(feature_test[0], target_test[0], color=train_color, label="train")
60
61
62
63
64 ### draw the regression line, once it's coded
65 try:
66     plt.plot( feature_test, reg.predict(feature_test) )
67 except NameError:
68     pass
69 plt.xlabel(features_list[1])
70 plt.ylabel(features_list[0])
71 plt.legend()
72 plt.show()

```

**File: ./gitignore**

```

*.pyc
enron_mail_20110402.tgz
enron_mail_20110402/
enron_mail_20150507.tgz
maildir/
text_learning/your_word_data.pkl
text_learning/your_email_authors.pkl
my_classifier.pkl
my_dataset.pkl
my_feature_list.pkl
.idea
.ipynb_checkpoints

```

**File: ./evaluation/evaluate\_poi\_identifier.py**

```

1  #!/usr/bin/python
2
3
4  """
5      Starter code for the evaluation mini-project.
6      Start by copying your trained/tested POI identifier from
7      that which you built in the validation mini-project.
8
9      This is the second step toward building your POI identifier!
10
11     Start by loading/formatting the data...
12 """
13
14 import pickle
15 import sys
16 sys.path.append("../tools/")
17 from feature_format import featureFormat, targetFeatureSplit
18
19 data_dict = pickle.load(open("../final_project/final_project_dataset.pkl", "rb" ) )
20
21 ### add more features to features_list!
22 features_list = ["poi", "salary"]
23
24 data = featureFormat(data_dict, features_list)
25 labels, features = targetFeatureSplit(data)
26
27
28
29 ### your code goes here
30
31

```

## File: ./k\_means/k\_means\_cluster.py

```
1 #!/usr/bin/python
2
3 """
4     Skeleton code for k-means clustering mini-project.
5 """
6
7
8
9
10 import pickle
11 import numpy
12 import matplotlib.pyplot as plt
13 import sys
14 sys.path.append("../tools/")
15 from feature_format import featureFormat, targetFeatureSplit
16
17
18
19
20 def Draw(pred, features, poi, mark_poi=False, name="image.png", f1_name="feature 1", f2_name="feature 2"):
21     """ some plotting code designed to help you visualize your clusters """
22
23     ### plot each cluster with a different color--add more colors for
24     ### drawing more than five clusters
25     colors = ["b", "c", "k", "m", "g"]
26     for ii, pp in enumerate(pred):
27         plt.scatter(features[ii][0], features[ii][1], color = colors[pred[ii]])
28
29     ### if you like, place red stars over points that are POIs (just for funsies)
30     if mark_poi:
31         for ii, pp in enumerate(pred):
32             if poi[ii]:
33                 plt.scatter(features[ii][0], features[ii][1], color="r", marker="*")
34     plt.xlabel(f1_name)
35     plt.ylabel(f2_name)
36     plt.savefig(name)
37     plt.show()
38
39
40
41 ### load in the dict of dicts containing all the data on each person in the dataset
42 data_dict = pickle.load(open("../final_project/final_project_dataset.pkl", "rb"))
43 ### there's an outlier--remove it!
44 data_dict.pop("TOTAL", 0)
45
46
47 ### the input features we want to use
48 ### can be any key in the person-level dictionary (salary, director_fees, etc.)
49 feature_1 = "salary"
50 feature_2 = "exercised_stock_options"
51 poi = "poi"
52 features_list = [poi, feature_1, feature_2]
53 data = featureFormat(data_dict, features_list)
54 poi, finance_features = targetFeatureSplit(data)
55
56
57 ### in the "clustering with 3 features" part of the mini-project,
58 ### you'll want to change this line to
59 ### for f1, f2, _ in finance_features:
60 ### (as it's currently written, the line below assumes 2 features)
61 for f1, f2 in finance_features:
62     plt.scatter(f1, f2)
63 plt.show()
64
65 ### cluster here; create predictions of the cluster labels
66 ### for the data and store them to a list called pred
67
68
69
70
71 ### rename the "name" parameter when you change the number of features
72 ### so that the figure gets saved to a different file
73 try:
74     Draw(pred, finance_features, poi, mark_poi=False, name="clusters.pdf", f1_name=feature_1, f2_name=feature_2)
75 except NameError:
76     print("no predictions object named pred found, no clusters to plot")
```

### File: ./choose\_your\_own/your\_algorithm.py

```
1 #!/usr/bin/python
2
3 import matplotlib.pyplot as plt
4 from prep_terrain_data import makeTerrainData
5 from class_vis import prettyPicture
6
7 features_train, labels_train, features_test, labels_test = makeTerrainData()
8
9
10 #### the training data (features_train, labels_train) have both "fast" and "slow"
11 #### points mixed together--separate them so we can give them different colors
12 #### in the scatterplot and identify them visually
13 grade_fast = [features_train[ii][0] for ii in range(0, len(features_train)) if labels_train[ii]==0]
14 bumpy_fast = [features_train[ii][1] for ii in range(0, len(features_train)) if labels_train[ii]==0]
15 grade_slow = [features_train[ii][0] for ii in range(0, len(features_train)) if labels_train[ii]==1]
16 bumpy_slow = [features_train[ii][1] for ii in range(0, len(features_train)) if labels_train[ii]==1]
17
18
19 ##### initial visualization
20 plt.xlim(0.0, 1.0)
21 plt.ylim(0.0, 1.0)
22 plt.scatter(bumpy_fast, grade_fast, color = "b", label="fast")
23 plt.scatter(grade_slow, bumpy_slow, color = "r", label="slow")
24 plt.legend()
25 plt.xlabel("bumpiness")
26 plt.ylabel("grade")
27 plt.show()
28 #####
29
30
31 #### your code here! name your classifier object clf if you want the
32 #### visualization code (prettyPicture) to show you the decision boundary
33
34
35
36
37
38
39
40
41 try:
42     prettyPicture(clf, features_test, labels_test)
43 except NameError:
44     pass
```

### File: ./choose\_your\_own/prep\_terrain\_data.py

```
1 #!/usr/bin/python
2 import random
3
4
5 def makeTerrainData(n_points=1000):
6     #####
7     #### make the toy dataset
8     random.seed(42)
9     grade = [random.random() for ii in range(0,n_points)]
10    bumpy = [random.random() for ii in range(0,n_points)]
11    error = [random.random() for ii in range(0,n_points)]
12    y = [round(grade[ii]*bumpy[ii]+0.3+0.1*error[ii]) for ii in range(0,n_points)]
13    for ii in range(0, len(y)):
14        if grade[ii]>0.8 or bumpy[ii]>0.8:
15            y[ii] = 1.0
16
17    #### split into train/test sets
18    X = [(gg, ss) for gg, ss in zip(grade, bumpy)]
19    split = int(0.75*n_points)
20    X_train = X[:split]
21    X_test = X[split:]
22    y_train = y[:split]
23    y_test = y[split:]
24
25    grade_sig = [X_train[ii][0] for ii in range(0, len(X_train)) if y_train[ii]==0]
26    bumpy_sig = [X_train[ii][1] for ii in range(0, len(X_train)) if y_train[ii]==0]
```

```

27 grade_bkg = [X_train[ii][0] for ii in range(0, len(X_train)) if y_train[ii]==1]
28 bumpy_bkg = [X_train[ii][1] for ii in range(0, len(X_train)) if y_train[ii]==1]
29
30 training_data = {"fast": {"grade": grade_sig, "bumpiness": bumpy_sig}
31                  , "slow": {"grade": grade_bkg, "bumpiness": bumpy_bkg}}
32
33
34 grade_sig = [X_test[ii][0] for ii in range(0, len(X_test)) if y_test[ii]==0]
35 bumpy_sig = [X_test[ii][1] for ii in range(0, len(X_test)) if y_test[ii]==0]
36 grade_bkg = [X_test[ii][0] for ii in range(0, len(X_test)) if y_test[ii]==1]
37 bumpy_bkg = [X_test[ii][1] for ii in range(0, len(X_test)) if y_test[ii]==1]
38
39 test_data = {"fast": {"grade": grade_sig, "bumpiness": bumpy_sig}
40             , "slow": {"grade": grade_bkg, "bumpiness": bumpy_bkg}}
41
42 return X_train, y_train, X_test, y_test
43

```

## File: ./choose\_your\_own/class\_vis.py

```

1  #!/usr/bin/python
2
3  import numpy as np
4  import matplotlib.pyplot as plt
5  import pylab as pl
6
7  def prettyPicture(clf, X_test, y_test):
8      x_min = 0.0; x_max = 1.0
9      y_min = 0.0; y_max = 1.0
10
11     # Plot the decision boundary. For that, we will assign a color to each
12     # point in the mesh [x_min, m_max]x[y_min, y_max].
13     h = .01 # step size in the mesh
14     xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
15     Z = clf.predict(np.c_[xx.ravel(), yy.ravel()])
16
17     # Put the result into a color plot
18     Z = Z.reshape(xx.shape)
19     plt.xlim(xx.min(), xx.max())
20     plt.ylim(yy.min(), yy.max())
21
22     plt.pcolormesh(xx, yy, Z, cmap=plt.cm.seismic)
23
24     # Plot also the test points
25     grade_sig = [X_test[ii][0] for ii in range(0, len(X_test)) if y_test[ii]==0]
26     bumpy_sig = [X_test[ii][1] for ii in range(0, len(X_test)) if y_test[ii]==0]
27     grade_bkg = [X_test[ii][0] for ii in range(0, len(X_test)) if y_test[ii]==1]
28     bumpy_bkg = [X_test[ii][1] for ii in range(0, len(X_test)) if y_test[ii]==1]
29
30     plt.scatter(grade_sig, bumpy_sig, color = "b", label="fast")
31     plt.scatter(grade_bkg, bumpy_bkg, color = "r", label="slow")
32     plt.legend()
33     plt.xlabel("bumpiness")
34     plt.ylabel("grade")
35
36     plt.savefig("test.png")
37
38 import base64
39 import json
40 import subprocess
41
42 def output_image(name, format, bytes):
43     image_start = "BEGIN_IMAGE_f9825uweof8jw9fj4r8"
44     image_end = "END_IMAGE_0238jfw08fjsiufhw8frs"
45     data = {}
46     data['name'] = name
47     data['format'] = format
48     data['bytes'] = base64.encodestring(bytes)
49     print (image_start+json.dumps(data)+image_end)
50

```