

```

1 import collections
2
3 import numpy as np
4
5 import util
6 import svm
7
8
9 def get_words(message):
10     """Get the normalized list of words from a message string.
11
12     This function should split a message into words, normalize them, and return
13     the resulting list. For splitting, you should split on spaces. For normalization,
14     you should convert everything to lowercase.
15
16     Args:
17         message: A string containing an SMS message
18
19     Returns:
20         The list of normalized words from the message.
21     """
22
23     # *** START CODE HERE ***
24     return [word.lower() for word in message.split(' ')]
25     # *** END CODE HERE ***
26
27
28 def create_dictionary(messages):
29     """Create a dictionary mapping words to integer indices.
30
31     This function should create a dictionary of word to indices using the provided
32     training messages. Use get_words to process each message.
33
34     Rare words are often not useful for modeling. Please only add words to the dictionary
35     if they occur in at least five messages.
36
37     Args:
38         messages: A list of strings containing SMS messages
39
40     Returns:
41         A python dict mapping words to integers.
42     """
43
44     # *** START CODE HERE ***
45     word_counts = collections.defaultdict(int)
46
47     for message in messages:
48         for word in set(get_words(message)):
49             word_counts[word] += 1
50
51     resulting_dictionary = {}
52
53     for word, count in word_counts.items():
54         if count >= 5:
55             next_index = len(resulting_dictionary)
56             resulting_dictionary[word] = next_index
57
58     print("RESULTING DICTIONARY: ", resulting_dictionary)
59     return resulting_dictionary
60     # *** END CODE HERE ***
61
62
63 def transform_text(messages, word_dictionary):
64     """Transform a list of text messages into a numpy array for further processing.
65
66     This function should create a numpy array that contains the number of times each word
67     of the vocabulary appears in each message.
68     Each row in the resulting array should correspond to each message
69     and each column should correspond to a word of the vocabulary.
70
71     Use the provided word dictionary to map words to column indices. Ignore words that
72     are not present in the dictionary. Use get_words to get the words for a message.

```

```

73
74 Args:
75     messages: A list of strings where each string is an SMS message.
76     word_dictionary: A python dict mapping words to integers.
77
78 Returns:
79     A numpy array marking the words present in each message.
80     Where the component (i,j) is the number of occurrences of the
81     j-th vocabulary word in the i-th message.
82 """
83 # *** START CODE HERE ***
84 result = np.zeros((len(messages), len(word_dictionary)))
85
86 for i, message in enumerate(messages):
87     for word in get_words(message):
88         if word in word_dictionary:
89             result[i, word_dictionary[word]] += 1
90
91 return result
92 # *** END CODE HERE ***
93
94
95 def fit_naive_bayes_model(matrix, labels):
96     """Fit a naive bayes model.
97
98     This function should fit a Naive Bayes model given a training matrix and labels.
99
100    The function should return the state of that model.
101
102    Feel free to use whatever datatype you wish for the state of the model.
103
104    Args:
105        matrix: A numpy array containing word counts for the training data
106        labels: The binary (0 or 1) labels for that training data
107
108    Returns: The trained model
109    """
110
111    # *** START CODE HERE ***
112    model = {}
113
114    phi = 1. * sum(labels) / len(labels)
115    model['logphi_0'] = np.log(1.-phi)
116    model['logphi_1'] = np.log(phi)
117    theta_0 = (matrix[labels == 0]).sum(axis=0) + 1
118    theta_1 = (matrix[labels == 1]).sum(axis=0) + 1
119    theta_0 /= theta_0.sum()
120    theta_1 /= theta_1.sum()
121    model['logtheta_0'] = np.log(theta_0)
122    model['logtheta_1'] = np.log(theta_1)
123
124    return model
125    # *** END CODE HERE ***
126
127
128 def predict_from_naive_bayes_model(model, matrix):
129     """Use a Naive Bayes model to compute predictions for a target matrix.
130
131     This function should be able to predict on the models that fit_naive_bayes_model
132     outputs.
133
134    Args:
135        model: A trained model from fit_naive_bayes_model
136        matrix: A numpy array containing word counts
137
138    Returns: A numpy array containg the predictions from the model
139    """
140    # *** START CODE HERE ***
141    output = np.zeros(matrix.shape[0])
142
143    logphi_0 = model['logphi_0']
144    logphi_1 = model['logphi_1']

```



```

145 logtheta_0 = model['logtheta_0']
146 logtheta_1 = model['logtheta_1']
147 logprobs_0 = (matrix * logtheta_0).sum(axis=1) + logphi_0
148 logprobs_1 = (matrix * logtheta_1).sum(axis=1) + logphi_1
149
150 output = (logprobs_1 > logprobs_0).astype(int)
151 return output
152 # *** END CODE HERE ***
153
154
155 def get_top_five_naive_bayes_words(model, dictionary):
156     """Compute the top five words that are most indicative of the spam (i.e positive) class.
157
158     Uses the metric given in part-c as a measure of how indicative a word is.
159     Return the words in sorted form, with the most indicative word first.
160
161     Args:
162         model: The Naive Bayes model returned from fit_naive_bayes_model
163         dictionary: A mapping of word to integer ids
164
165     Returns: A list of the top five most indicative words in sorted order with the most indicative first
166     """
167     # *** START CODE HERE ***
168     ids = np.argsort(model['logtheta_0'] - model['logtheta_1'])[:5]
169
170     reverse_dictionary = {i: word for word, i in dictionary.items()}
171
172     return [reverse_dictionary[i] for i in ids]
173     # *** END CODE HERE ***
174
175
176 def compute_best_svm_radius(train_matrix, train_labels, val_matrix, val_labels, radius_to_consider):
177     """Compute the optimal SVM radius using the provided training and evaluation datasets.
178
179     You should only consider radius values within the radius_to_consider list.
180     You should use accuracy as a metric for comparing the different radius values.
181
182     Args:
183         train_matrix: The word counts for the training data
184         train_labels: The spma or not spam labels for the training data
185         val_matrix: The word counts for the validation data
186         val_labels: The spam or not spam labels for the validation data
187         radius_to_consider: The radius values to consider
188
189     Returns:
190         The best radius which maximizes SVM accuracy.
191     """
192     # *** START CODE HERE ***
193
194     best = None
195
196     for radius in radius_to_consider:
197         svm_predictions = svm.train_and_predict_svm(train_matrix, train_labels, val_matrix, radius)
198         svm_accuracy = np.mean(svm_predictions == val_labels)
199
200         if best is None:
201             best = (svm_accuracy, radius)
202         else:
203             best = max(best, (svm_accuracy, radius))
204
205     return best[1]
206     # *** END CODE HERE ***
207
208
209 def main():
210     train_messages, train_labels = util.load_spam_dataset('spam_train.tsv')
211     val_messages, val_labels = util.load_spam_dataset('spam_val.tsv')
212     test_messages, test_labels = util.load_spam_dataset('spam_test.tsv')
213
214     dictionary = create_dictionary(train_messages)
215
216     print('Size of dictionary: ', len(dictionary))

```

```

217 util.write_json('spam_dictionary', dictionary)
218
219 train_matrix = transform_text(train_messages, dictionary)
220
221 np.savetxt('spam_sample_train_matrix', train_matrix[:100,:])
222
223 val_matrix = transform_text(val_messages, dictionary)
224 test_matrix = transform_text(test_messages, dictionary)
225
226 naive_bayes_model = fit_naive_bayes_model(train_matrix, train_labels)
227
228 naive_bayes_predictions = predict_from_naive_bayes_model(naive_bayes_model, test_matrix)
229
230 np.savetxt('spam_naive_bayes_predictions', naive_bayes_predictions)
231
232 naive_bayes_accuracy = np.mean(naive_bayes_predictions == test_labels)
233
234 print('Naive Bayes had an accuracy of {} on the testing set'.format(naive_bayes_accuracy))
235
236 top_5_words = get_top_five_naive_bayes_words(naive_bayes_model, dictionary)
237
238 print('The top 5 indicative words for Naive Bayes are: ', top_5_words)
239
240 util.write_json('spam_top_indicative_words', top_5_words)
241
242 optimal_radius = compute_best_svm_radius(train_matrix, train_labels, val_matrix, val_labels, [0.01, 0.1, 1, 10])
243
244 util.write_json('spam_optimal_radius', optimal_radius)
245
246 print('The optimal SVM radius was {}'.format(optimal_radius))
247
248 svm_predictions = svm.train_and_predict_svm(train_matrix, train_labels, test_matrix, optimal_radius)
249
250 svm_accuracy = np.mean(svm_predictions == test_labels)
251
252 print('The SVM model had an accuracy of {} on the testing set'.format(svm_accuracy, optimal_radius))
253
254
255
256 if __name__ == "__main__":
257     main()

```