

```

1
2 # CIS 511 NLP - Assignment 2 - Viterbi Part-of-speech Tagger
3
4 """
5 Created on Sat Feb 25 08:06:51 2020
6
7 @author: Siyu Yang
8 @unique name: siyuya
9 @UMID:76998080
10 """
11
12 import numpy as np
13 import sys
14
15 #=== Functions ===#
16
17 # 1) Load Train File
18 def load_train_data(trainFile):
19
20     # collect and store raw counts required for algorithm
21
22     tagFreqs_dict = {} # store corresponding frequency for each unique tag (KEY:
23     tag, VALUE: count)
24     wordFreqs_dict = {} # store corresponding frequency for each unique word (KEY:
25     word, VALUE: count)
26     tag_index = 0
27     word_index = 0
28
29     train_sample = []
30
31     with open(trainFile) as data:
32         for line in data:
33             line = "<s>/<s>" + ' ' + line # process files as sentences
34             splitline = line.strip().split() # add beginning
35
36             new_line = []
37
38             for pair in splitline:
39                 pair = pair.split('/')
40                 if len(pair) == 2:
41                     new_line.append(pair)
42
43             train_sample.append(new_line)
44
45             for i in range(len(new_line)):
46                 pre_tag = new_line[i][1]
47
48                 if pre_tag not in tagFreqs_dict: # add unique tag into dict
49                     tagFreqs_dict[pre_tag] = tag_index
50                     tag_index += 1
51
52                 word = new_line[i][0]
53
54                 if word not in wordFreqs_dict: # add unique word into dict
55                     wordFreqs_dict[word] = word_index
56                     word_index += 1
57
58             wordFreqs_dict['NON'] = word_index
59         return tagFreqs_dict, wordFreqs_dict, train_sample
60
61 # 1) Process Train File
62
63 def process_train_data(train_sample, tagFreqs_dict, wordFreqs_dict):
64
65     tag_num = len(tagFreqs_dict) # total number of tags
66     word_num = len(wordFreqs_dict) # total number of words
67
68     tag_to_tag = np.ones((tag_num, tag_num))
69     tag_to_word = np.ones((tag_num, word_num))
70

```

```

71
72
73     for new_line in train_sample:
74         for i in range(len(new_line)):
75             word = new_line[i][0]
76             pre_tag = new_line[i][1]
77             tag_to_word[tagFreqs_dict[pre_tag],wordFreqs_dict[word]] += 1
78             if i+1 < len(new_line):
79                 next_tag = new_line[i+1][1]
80                 tag_to_tag[tagFreqs_dict[pre_tag], tagFreqs_dict[next_tag]] += 1
81
82     tag_to_word = (tag_to_word/tag_to_word.sum(axis=1, keepdims=1))
83     tag_to_tag = tag_to_tag/tag_to_tag.sum(axis=1, keepdims=1)
84
85     return tag_to_word, tag_to_tag
86
87
88 # 3) Load Test File
89
90 def load_test_data(testFile):
91     test_sample = []
92     with open(testFile) as data:
93         for line in data:
94             line = "<s>/<s>" + ' ' + line
95             splitline = line.strip().split()
96             new_line = []
97             for pair in splitline:
98                 pair = pair.split('/')
99                 if len(pair) == 2:
100                     new_line.append(pair)
101             test_sample.append(new_line)
102     return test_sample
103
104
105
106 # 4) Viterbi Algorithm
107
108 def Viterbi(test_sample,tag_to_word,tag_to_tag, tagFreqs_dict, wordFreqs_dict):
109     # create storage for sentence tag predictions and true
110     sen_samples_predit = []
111     sen_samples_true = []
112
113     # iterate through every test sentence
114     for instance in test_sample:
115
116         # 0) create storage for this sentence's predictions, probability scores, and
117         # back tag pointer
118
119         prob_score = np.zeros((len(tagFreqs_dict),len(instance)))
120         back_tag = np.zeros((len(tagFreqs_dict),len(instance)))
121         label_list = np.zeros(len(instance))
122
123         # 1) Initialization
124         tag_init = '<s>'
125         word = instance[1][0]
126
127         if word not in wordFreqs_dict:
128             word = 'NON'
129
130         label_list[1] = (tagFreqs_dict[instance[1][1]])
131
132         word_given_tag_prob = tag_to_word[:, wordFreqs_dict[word]]
133         tag_given_tag_prob = tag_to_tag[0,:]
134
135         prob_score[:,1] = list(word_given_tag_prob * tag_given_tag_prob) # list
136         back_tag[:,1] = np.zeros((len(tagFreqs_dict))) # index
137
138         # 2) Iteration
139         for i in range(2,len(instance)):
140             word = instance[i][0]
141
142             max_score = [] # max score for each one

```

```

142         max_back_tag = [] # max score position
143
144         if word not in wordFreqs_dict:
145             word = 'NON'
146         label_list[i] = (tagFreqs_dict[instance[i][1]])
147
148         transition_score = np.array(prob_score[:, i-1]).reshape(-1, 1) *
tag_to_tag
149
150         for j in range(len(tagFreqs_dict)): # for each tag
151
152             word_pre_prob = tag_to_word[j, wordFreqs_dict[word]]
153
154             max_score.append(max(transition_score[:, j] * word_pre_prob)) #
score for this tag
155             max_back_tag.append(np.argmax(transition_score[:, j] * word_pre_prob))
156
157         prob_score[:, i] = np.array(max_score)
158         back_tag[:, i] = np.array(max_back_tag)
159
160         # 3) store predictions for this sentence
161
162         final_max = np.argmax(prob_score[:, -1])
163
164         predict_label = np.zeros((len(instance)))
165
166         predict_label[-1] = int(final_max)
167
168         for i in range(len(instance)-2, 0, -1):
169             final_max = int(back_tag[final_max, i+1])
170             predict_label[i] = int(final_max)
171
172
173         sen_samples_predit.append(predict_label)
174         sen_samples_true.append(label_list)
175
176         return sen_samples_predit, sen_samples_true
177
178
179     # 5) Calculate Predict Score
180
181     def predict(sen_samples_predit, sen_samples_true):
182         # create storage for sentence tag predictions
183
184         sentence_num = 0
185         tag_num = 0
186         correct_sentence_num = 0
187         wrong_tag_num = 0
188         for sentence in range(len(sen_samples_predit)):
189             sentence_num += 1
190             flag = 1
191             for tag in range(len(sen_samples_predit[sentence])):
192                 tag_num += 1
193                 if (sen_samples_predit[sentence][tag]) != (sen_samples_true[sentence][tag]):
194                     flag = 0
195                     wrong_tag_num += 1
196             if flag == 1:
197                 correct_sentence_num += 1
198
199         tag_accuracy = (tag_num - wrong_tag_num) / tag_num
200         sen_accuracy = correct_sentence_num / sentence_num
201
202         return tag_accuracy, sen_accuracy
203
204
205     def ouput(sen_samples_predit, tagFreqs_dict):
206         index_tag_dict = {}
207         for tag, index in tagFreqs_dict.items():
208             index_tag_dict[index] = tag
209         data_path = sys.argv[2]
210

```

```

211 test_sample = []
212 index = 0
213 with open(data_path) as data:
214     for line in data:
215         splitline = line.strip().split()
216         new_line = ''
217         item_index = 0
218         for pair in splitline:
219             pair = pair.split('/')
220             word = pair[0]
221             if len(pair) == 2:
222                 tag = pair[1]
223                 new_line += word + '/' + index_tag_dict[sen_samples_predict[index]]
224                 item_index += 1
225             index += 1
226         test_sample.append(new_line)
227 with open('POS.test.out', 'w') as f:
228     for sentenc in test_sample:
229         f.write(sentenc)
230         f.write('\n')
231
232
233 # 6) A Simple Baseline Program
234
235 def baseline(trainFile, testFile, tag_to_word, tag_to_tag, tagFreqs_dict,
wordFreqs_dict):
236     tag_count = 0
237     accuracy_count = 0
238     tag_total = len(tagFreqs_dict)
239     word_total = len(wordFreqs_dict)
240     word_tag = np.ones((word_total, tag_total))
241
242     with open(trainFile) as data:
243         for line in data:
244             splitline = line.strip().split()
245             new_line = []
246             for pair in splitline:
247                 pair = pair.split('/')
248                 if len(pair) == 2:
249                     new_line.append(pair)
250             for i in range(len(new_line)):
251                 word = new_line[i][0]
252                 pre_tag = new_line[i][1]
253                 if word not in wordFreqs_dict:
254                     word = 'NON'
255                 word_tag[wordFreqs_dict[word], tagFreqs_dict[pre_tag]] += 1
256 word_tag = word_tag / word_tag.sum(axis=1, keepdims=1)
257
258     with open(testFile) as data:
259         for line in data:
260             splitline = line.strip().split()
261             for pair in splitline:
262                 pair = pair.split('/')
263                 if len(pair) == 2:
264                     tag_count += 1
265                     word = ''
266                     if pair[0] not in wordFreqs_dict:
267                         word = 'NON'
268                     else:
269                         word = pair[0]
270                     predict_tag = np.argmax(word_tag[wordFreqs_dict[word], :])
271                     if predict_tag == tagFreqs_dict[pair[1]]:
272                         accuracy_count += 1
273
274
275 baseline_accuracy = accuracy_count / tag_count
276 return baseline_accuracy
277
278
279
280

```

```
281
282
283
284 if __name__ == '__main__':
285     # 1) Load files
286     trainFile = sys.argv[1]
287     testFile = sys.argv[2]
288
289     # 2) Process train & test data
290     tagFreqs_dict, wordFreqs_dict, train_sample = load_train_data(trainFile)
291
292     tag_to_word, tag_to_tag = process_train_data(train_sample, tagFreqs_dict,
293                                                  wordFreqs_dict)
294
295     test_sample = load_test_data(testFile)
296
297     # 3) Run Viterbi algorithm
298     sen_samples_predit, sen_samples_true = Viterbi(test_sample, tag_to_word, tag_to_tag
299                                                    , tagFreqs_dict, wordFreqs_dict)
300
301     # 4) Calculate predict score
302     tag_accuracy, sen_accuracy = predict(sen_samples_predit, sen_samples_true)
303
304     # 5) Baseline
305     baseline_accuracy = baseline(trainFile, testFile, tag_to_word, tag_to_tag,
306                                  tagFreqs_dict, wordFreqs_dict)
307
308     ouput(sen_samples_predit, tagFreqs_dict)
309
310     # print accuracy
311     print ("Viterbi tag accuracy is ", tag_accuracy)
312     print ('Baseline accuracy is ', baseline_accuracy)
```