

```

1  # CIS 511 NLP - Assignment 4.2 - Dialog Act Classification
2
3  """
4  Created on Sat Apr 11 15:20:36 2020
5
6  @author: Siyu Yang
7  @unique name: siyuya
8  @UMID:76998080
9  """
10
11 import sys
12 import math
13
14 #=== Functions ===#
15
16 # 1) find the middle content between begin word and end word
17 def middle_content(begin, end, content):
18     mid_content = ''
19     if content.find(begin): # find the line start with begin words
20         beginword = content[content.find(begin):content.rfind(end)]
21         mid_content = beginword[len(begin):]
22     return mid_content
23
24
25 # 2) Parse Training File
26 def parse_file(filename):
27
28     # dialog_dict: key = DialogAct, value = list of all words
29     dialog_dict = dict()
30     # number_dict: key = DialogAct, value = times in train data
31     number_dict = dict()
32
33     with open(filename, 'r', encoding='utf-8') as f:
34         content = f.read()
35         lines = content.split("\n")
36         for line in lines:
37
38             if len(line) > 2:
39                 prev_words = []
40                 label_answers = []
41
42                 if line.startswith("Student:"):
43                     words = line.split(" ")
44                     for word in words:
45                         label_answers.append(word)
46
47                 prev_words = label_answers
48
49                 if line.startswith("Advisor:"):
50                     dialog_act = middle_content("[", "]", line)
51
52                     if dialog_act == "social" or dialog_act == "pull" or dialog_act
53                     == "push":
54                         dialog_act = ""
55                     else:
56                         if dialog_act not in number_dict:
57                             number_dict[dialog_act] = 1
58                         else:
59                             number_dict[dialog_act] += 1
60
61                 #student's words and advisor's words are both written to
62                 if len(prev_words) != 0 and len(dialog_act) != 0:
63                     if dialog_act in dialog_dict:
64                         for word in prev_words:
65                             dialog_dict[dialog_act].append(word)
66                     else:
67                         dialog_dict[dialog_act] = prev_words
68
69                 prev_words = []
70                 dialog_act = ""
71
72     return dialog_dict, number_dict

```

```

72
73
74 # 3) Generate Probability Dictionary
75 def probability(number_dict):
76     # prob_dict: key = sense , value = probability of sense
77     prob_dict = dict()
78
79     total = 0
80     for sense in number_dict:
81         total += number_dict[sense]
82     for sense in number_dict:
83         prob_dict[sense] = number_dict[sense] / total
84
85     return prob_dict
86
87
88 # 4) Generate Unique Dictionary
89 def generate_unique_dict(dialog_dict):
90     # unique_dict: key = DialogAct, value = list of unique words
91     unique_dict = dict()
92
93     for dialog_act in dialog_dict:
94         unique_dict[dialog_act] = list()
95         for word in dialog_dict[dialog_act]:
96             if not word in unique_dict[dialog_act]:
97                 unique_dict[dialog_act].append(word)
98
99     return unique_dict
100
101
102 # 5) Process Test Data
103 def process_test_data(filename):
104     # test_dict: key = ID of the line , value = list of all words
105     test_dict = dict()
106
107     ID = 0
108     with open(filename, 'r', encoding='utf-8') as f:
109         content = f.read()
110         lines = content.split("\n")
111         for line in lines:
112             if line.startswith("Student:"):
113                 words = line.split(" ")
114                 for word in words:
115                     if word != "Student:":
116                         if ID in test_dict:
117                             test_dict[ID].append(word)
118                         else:
119                             test_dict[ID] = [word]
120                 ID += 1
121     return test_dict
122
123
124 # 6) Add One Smoothing and Get Output
125 def add_one_smoothing(test_dict, dialog_dict, number_dict, unique_dict, prob_dict,
126 output):
127     # score_dict: key = ID , value = sense and score
128     score_dict = dict()
129     # final_dict: key = ID , value = final label
130     final_dict = dict()
131
132     for ID in test_dict:
133         for sense in dialog_dict:
134             total = 0
135             for word in test_dict[ID]:
136                 # number of word appearance in sense
137                 num_1 = dialog_dict[sense].count(word) + 1
138                 # number of sense appearance
139                 num_2 = number_dict[sense] + len(unique_dict[sense])
140
141                 total *= math.log((num_1 / num_2), 2)
142

```

```

143         score = math.log(prob_dict[sense], 2) + total
144
145         if ID not in score_dict:
146             score_dict[ID] = list()
147             score_dict[ID].append({sense:score})
148
149
150     # iterates a dictionary to find argmax and best sense for ID
151     final_label = ""
152     argMax = -999
153     for ID in score_dict:
154         for scores in score_dict[ID]:
155             for sense in scores:
156                 score = scores[sense]
157                 if score > argMax:
158                     argMax = score
159                     final_label = sense
160             final_dict[ID] = final_label
161
162     for ID in final_dict:
163         for word in test_dict[ID]:
164             output.write(word+" ")
165             output.write("\n"+"Label: "+str(final_dict[ID])+"\n\n")
166
167     return final_dict
168
169
170 # 7) Calculate The Accuracy
171 def cal_acc(final_dict, testfile):
172
173     # test_out_dict: key = ID , value = test data
174     test_out_dict = dict()
175
176     with open(testfile, encoding='utf-8' ) as f:
177         line_num = 0
178         content = f.read()
179         lines = content.split("\n")
180         for line in lines:
181             if len(line) > 2:
182                 if line.startswith("Advisor:"):
183                     dialog_act = middle_content("[", "]", line)
184                 if line.startswith("Student:"):
185                     if len(dialog_act) > 0:
186                         test_out_dict[line_num] = dialog_act
187                     line_num += 1
188
189
190     correct_num = 0
191     total_num = 0
192     for word in final_dict:
193         if word in final_dict and word in test_out_dict:
194             if final_dict[word] == test_out_dict[word]:
195                 correct_num += 1
196             total_num += 1
197     accuracy = correct_num/total_num
198     print("Accuracy:", accuracy)
199
200     return accuracy
201
202
203
204
205
206 if __name__ == '__main__':
207
208     # 1) load file
209     trainFile = sys.argv[1]
210     testFile = sys.argv[2]
211
212     # 2) create outputfile
213     outputName = testFile + ".out"
214     output = open(outputName, "w")

```

```
215
216 # 3) Parse Training File
217 dialog_dict, number_dict = parse_file(trainFile)
218
219 # 4) Generate Probability Dictionary
220 prob_dict = probability(number_dict)
221
222 # 5) Generate Unique Dictionary
223 unique_dict = generate_unique_dict(dialog_dict)
224
225 # 6) Process Test Data
226 test_dict = process_test_data(testFile)
227
228 # 7) Add One Smoothing and Get Output
229 final_dict = add_one_smoothing(test_dict, dialog_dict, number_dict, unique_dict,
230                                prob_dict, output)
231
232 # 8) Calculate The Accuracy
233 cal_acc(final_dict, testFile)
```