```python
# CIS 511 NLP - Assignment 4.1 - Natural Language Understanding for Dialog Systems

"""
Created on Sat April 11 20:30:18 2020

@author: Siyu Yang
@unique name: siyuya
@UMID:76998080
"""

from collections import defaultdict
import pandas
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score
from sklearn.model_selection import train_test_split
import re
import sys


#=== Functions ===#

#   1)   find the middle content between begin word and end word
def middle_content(begin, end, content):
    mid_content  =''

    if content.find(begin): # find the line start with begin words
        beginword = content[content.find(begin):content.rfind(end)]
        mid_content = beginword[len(begin):]
        return mid_content


#   2)   convert I/O/B into number 1/0/2
def IOB_to_Num(x):
    IOB_Char = x[-1]
    if IOB_Char=='I':
        return 1
    elif IOB_Char =='B':
        return 2
    else :
        return 0



#   3)   Process Train File to get train data
def Process_TrainFile(trainfile):

    trainList = open(trainfile,'r').read() #read the file and remove the empty lines.
    name_list = []
    train_fea_dic = defaultdict(list) # use default dictionary as train file feature
    dictionary


    for line in trainList.split("\n\n"):

        if("<class" in line):
            text = line.split("\n")
            text1 = text[0]
            text2 = re.sub('([.,!?()])', r' \1 ', text1)
            text3 = re.sub('\s{2,}', ' ', text2)
            value= text3.split(" ") #get the value from each line

            while "" in value: #remove empty lines
                value.remove("")

            for i in range(0, len(value)):
                # Feature 1:  the value of the token
                train_fea_dic['Value'].append(value[i])

                # Feature 2:  is token all uppercase?
                if(value[i].isupper):
                    train_fea_dic['UpperCase'].append(1)
                else:
```

```python
 72                                  train_fea_dic['UpperCase'].append(0)
 73
 74                          # Feature 3:  does token start with capital?
 75                          if(value[i].istitle()):
 76                              train_fea_dic['Capital'].append(1)
 77                          else:
 78                              train_fea_dic['Capital'].append(0)
 79
 80                          # Feature 4:  length of token
 81                          train_fea_dic['Length'].append(len(value[i]))
 82
 83                          # Feature 5:  does the token consist only of numbers?
 84                          if(value[i].isnumeric()):
 85                              train_fea_dic['Numeric'].append(1)
 86                          else:
 87                              train_fea_dic['Numeric'].append(0)
 88
 89                          # Feature 6:  does token start with vowel?
 90                          if value[i] in 'aeiou':
 91                              train_fea_dic['Vowel'].append(1)
 92                          else:
 93                              train_fea_dic['Vowel'].append(0)
 94
 95                          # Feature 7:  length of left word < 4
 96                          if(len(value[i])<4):
 97                              train_fea_dic['Length<Four'].append(1)
 98                          else:
 99                              train_fea_dic['Length<Four'].append(0)
100
101                          # Feature 8:  if Right word contains "."
102                          if('.' in value[i]):
103                              train_fea_dic['Period_in_token'].append(1)
104                          else:
105                              train_fea_dic['Period_in_token'].append(0)
106
107
108
109
110                  if text[1].startswith("<class"):
111                      mid_content = middle_content("<class", ">", line)
112                      content = mid_content.split("\n")
113
114                      for num in range(0,len(content)):
115                          if("id=" in content[num]):
116                              ID_value = content[num].split("=")[1]
117                          if("name=" in content[num]):
118                              name_list = content[num].split("=")[1].split(" ")
119
120
121                  for i in range(0, len(value)):
122                      for j in range(0,len(name_list)):
123                          if(value[i] == name_list[j]):
124                              if(j==0):
125                                  value[i]=name_list[j]+("/B")
126                              elif(j>0):
127                                  value[i]=name_list[j]+("/I")
128                      if(value[i] == ID_value):
129                          value[i]=value[i]+("/B")
130
131                      else:
132                          value[i]=value[i]+("/O")
133                      train_fea_dic['IOB'].append(value[i])
134
135
136          train_data = pandas.DataFrame.from_dict(train_fea_dic)
137          train_data['Value'] = train_data.index
138          train_data['IOB'] = train_data['IOB'].map(IOB_to_Num)
139
140          return train_data
141
142
143  #   4)  Process Test File to get test data
```

```python
144    def Process_TestFile(testfile):
145        testList = open(testfile,'r').read() #read the file and remove the empty lines.
146        name_list = []
147        test_fea_dic = defaultdict(list) # use default dictionary as test file feature
           dictionary
148
149        for line in testList.split("\n\n"):
150
151            if("<class" in line):
152                testtext = line.split("\n")
153                testtext1 = testtext[0]
154                testtext2 = re.sub('([.,!?()])', r' \1 ', testtext1)
155                testtext3 = re.sub('\s{2,}', ' ', testtext2)
156                testvalue= testtext3.split(" ") #get the value from each line
157
158
159                while "" in testvalue: #remove empty lines
160                    testvalue.remove("")
161
162                for i in range(0, len(testvalue)):
163                    # Feature 1:  the value of the token
164                    test_fea_dic['Value'].append(testvalue[i])
165
166                    # Feature 2:  is token all uppercase?
167                    if(testvalue[i].isupper):
168                        test_fea_dic['UpperCase'].append(1)
169                    else:
170                        test_fea_dic['UpperCase'].append(0)
171
172                    # Feature 3:  does token start with capital?
173                    if(testvalue[i].istitle()):
174                        test_fea_dic['Capital'].append(1)
175                    else:
176                        test_fea_dic['Capital'].append(0)
177
178                    # Feature 4:  length of token
179                    test_fea_dic['Length'].append(len(testvalue[i]))
180
181                    # Feature 5:  does the token consist only of numbers?
182                    if(testvalue[i].isnumeric()):
183                        test_fea_dic['Numeric'].append(1)
184                    else:
185                        test_fea_dic['Numeric'].append(0)
186
187                    # Feature 6:  does token start with vowel?
188                    if testvalue[i] in 'aeiou':
189                        test_fea_dic['Vowel'].append(1)
190                    else:
191                        test_fea_dic['Vowel'].append(0)
192
193
194                    # Feature 7:  length of left word < 4
195                    if(len(testvalue[i])<4):
196                        test_fea_dic['Length<Four'].append(1)
197                    else:
198                        test_fea_dic['Length<Four'].append(0)
199
200
201                    # Feature 8:  if Right word contains "."
202                    if('.' in testvalue[i]):
203                        test_fea_dic['Period_in_token'].append(1)
204                    else:
205                        test_fea_dic['Period_in_token'].append(0)
206
207
208                if testtext[1].startswith("<class"):
209                    mid_content = middle_content("<class", ">", line)
210                    content =mid_content.split("\n")
211
212
213                    for num in range(0,len(content)):
214                        if("name=" in content[num]):
```

```python
                            name_list = content[num].split("=")[1].split(" ")

                        if("id=" in content[num]):
                            ID_value = content[num].split("=")[1]

                for i in range(0, len(testvalue)):
                    for j in range(0,len(name_list)):
                        if(testvalue[i] == name_list[j]):
                            if(j==0):
                                testvalue[i]=name_list[j]+("/B")
                            elif(j>0):
                                testvalue[i]=name_list[j]+("/I")
                    if(testvalue[i] == ID_value):
                        testvalue[i]=testvalue[i]+("/O")

                    else:
                        testvalue[i]=testvalue[i]+("/O")
                    test_fea_dic['IOB'].append(testvalue[i])


        test_data = pandas.DataFrame.from_dict(test_fea_dic)
        test_data['Value'] = test_data.index
        test_data['IOB'] = test_data['IOB'].map(IOB_to_Num )

        return test_data,test_fea_dic

    #   5)  Calculate The Accuracy and Generate Output
    def accuracy(train_data, test_data,test_fea_dic,output):

        features = ['Value','UpperCase','Capital','Length','Numeric','Vowel',
        'Length<Four','Period_in_token']
        X = train_data[features]
        Y = train_data['IOB']
        test_X = test_data[features]
        test_Y = test_data['IOB']
        X_train,  X_test, Y_train, Y_test = train_test_split(X, Y)

        # Create Decision Tree classifer object and get the predict value
        predict = DecisionTreeClassifier().fit(X_train,Y_train).predict(test_X)

        # Calculate and print the accuracy score
        accuracy = accuracy_score(test_Y,predict)
        print("Accuracy:",accuracy)


        for i in range(0,len(predict)):
            output.write('\n')
            if(predict[i] == 0):
                output.write(str(test_fea_dic['Value'][i]))
                output.write('/O')
            elif(predict[i] == 1):
                output.write(str(test_fea_dic['Value'][i]))
                output.write('/I')
            else:
                output.write(str(test_fea_dic['Value'][i]))
                output.write('/B')

        output.close()
        return accuracy




    if __name__ == "__main__":

        # 1) load file
        trainfile = sys.argv[1]
        testfile = sys.argv[2]

        # 2) create outputfile
        outputName = testfile + ".out"
        output = open(outputName, "w")
```

```python
        # 2)  Process Train File to get train data
        train_data= Process_TrainFile(trainfile)

        # 3)  Process Test File to get test data
        test_data,test_fea_dic = Process_TestFile(testfile)

        # 4) generate output and calculate accurancy
        accuracy(train_data, test_data,test_fea_dic, output)
```