

Artificial and Computational Intelligence (Assignment - 2)

Problem Statement

As part of the 2nd Assignment, we'll try and predict the Part of Speech (POS) tag for each word in a provided sentence.

You are required to build a model using Hidden Markov Models which would help you predict the POS tags for all words in an utterance.

What is a POS tag?

In corpus linguistics, part-of-speech tagging (POS tagging or PoS tagging or POST), also called grammatical tagging or word-category disambiguation, is the process of marking up a word in a text (corpus) as corresponding to a particular part of speech, based on both its definition and its context—i.e., its relationship with adjacent and related words in a phrase, sentence, or paragraph. A simplified form of this is commonly taught to school-age children, in the identification of words as nouns, verbs, adjectives, adverbs, etc.

Dataset

The dataset can be downloaded from https://drive.google.com/open?id=1345iaxqTImJN6mKGh_c1T5n2OWumpyYz (https://drive.google.com/open?id=1345iaxqTImJN6mKGh_c1T5n2OWumpyYz). You can access it only using your BITS IDs.

Dataset Description

Sample Tuple

b100-5507

Mr. NOUN
Podger NOUN
had VERB
thanked VERB
him PRON
gravely ADV
, .
and CONJ
now ADV
he PRON
made VERB
use NOUN
of ADP
the DET
advice NOUN
. .

Explanation

The first token "b100-5507" is just a key and acts like an identifier to indicate the beginning of a sentence. The other tokens have a (Word, POS Tag) pairing.

List of POS Tags are: .

ADJ
ADP
ADV
CONJ
DET
NOUN
NUM
PRON
PRT
VERB
X

Note

. is used to indicate special characters such as '.', ',', ''

X is used to indicate vocab not part of English Language mostly. Others are Standard POS tags.

Evaluation

We wish to evaluate based on

- coding practices being followed
- commenting to explain the code and logic behind doing something
- your understanding and explanation of data
- how good the model would perform on unseen data.

Train-Test Split

Let us use a 80-20 split of our data for training and evaluation purpose.

In [1]:

```
1 #Import libraries
2 import random
3 from collections import Counter, defaultdict, namedtuple, OrderedDict
4 from itertools import chain
5 from pomegranate import State, HiddenMarkovModel, DiscreteDistribution
```

In [7]:

```

1 #Read data
2
3 Sentence = namedtuple("Sentence", "words tags")
4
5 # Function to read the data file and tokenize into sentences .
6 def read_data(filename):
7     """Read tagged sentence data"""
8     with open(filename, 'r') as f:
9         sentence_lines = [l.split("\n") for l in f.read().split("\n\n")]
10        sentences = OrderedDict(((s[0], Sentence(*zip(*[l.strip().split("\t")
11            for l in s[1:]]))) for s in sentence_lines if s[0]))
12    return sentences, sentence_lines
13 sentences, data = read_data('data.txt')
14
15 # Store all the predefined tags. Below logic can be moved into file to make it
16 # since tag is not stored in file we are just statically storing it .
17
18 tags = ['ADJ', 'ADP', 'ADV', 'CONJ', 'DET', 'NOUN', 'NUM', 'PRON', 'PRT', 'VERB', 'X']
19 tagset = frozenset(tags)
20 sentences

```

```

NJ', 'VERB', 'NOUN', 'DET', 'NOUN', 'ADP', 'NOUN', 'CONJ', 'DET', 'NO
UN', 'ADP', 'NOUN', '.'))',
    ('b100-54920',
        Sentence(words=('Could', 'it', 'just', 'be', ',', 'Ther
esa', 'wondered', ',', 'that', 'Anne', 'had', 'understood', 'only',
'too', 'well', ',', 'and', 'that', 'George', 'all', 'along', 'was',
'extraordinary', 'only', 'in', 'the', 'degree', 'to', 'which', 'he',
'was', 'dull', '?', '?'), tags=('VERB', 'PRON', 'ADV', 'VERB', '.',
'NOUN', 'VERB', '.', 'ADP', 'NOUN', 'VERB', 'VERB', 'ADV', 'ADV', 'AD
V', '.', 'CONJ', 'ADP', 'NOUN', 'PRT', 'ADV', 'VERB', 'ADJ', 'ADV',
'ADP', 'DET', 'NOUN', 'ADP', 'DET', 'PRON', 'VERB', 'ADJ', '.',
'.'))),
    ('b100-32087',
        Sentence(words=('Since', 'Af', 'and', 'P', 'divides',
'Af', 'for', 'Af', ',', 'we', 'have', 'Af', '.'), tags=('ADP', 'NOU
N', 'CONJ', 'NOUN', 'VERB', 'NOUN', 'ADP', 'NOUN', '.', 'PRON', 'VER
B', 'NOUN', '.'))),
    ('b100-30402', Sentence(words=('1', '.'), tags=('NUM',
'.'))),
    ('b100-17066',

```

In [3]:

```

1  #Pre-process data (Whatever you feel might be required)
2
3  ''' 3. UNIQUE KEY FOR EACH SENTENCE '''
4  keys = tuple(sentences.keys())
5
6  ''' 4. CONSTRUCT A VOCABULARY SET OF ALL UNIQUE WORDS '''
7  wordset = frozenset(chain(*[s.words for s in sentences.values()]))
8
9  ''' 5. WORD AND TAG SEQUENCES AS A TUPLE SET OF TUPLES '''
10 word_sequences = tuple([sentences[key].words for key in keys])
11 tag_sequences = tuple([sentences[key].tags for key in keys])
12
13 # written clas and function to split the training and test data which eventually
14 # it also has more frequent use logic like tag_Sequcning ,word_Sequencing, word
15
16 class Subset(namedtuple("BaseSet", "sentences keys vocab X tagset Y")):
17
18     def __new__(cls, sentences, keys):
19
20         word_sequences = tuple([sentences[key].words for key in keys])
21         tag_sequences = tuple([sentences[key].tags for key in keys])
22         wordset = frozenset(chain(*word_sequences))
23         tagset = frozenset(chain(*tag_sequences))
24
25         return super().__new__(cls,
26                                {key: sentences[key] for key in keys},
27                                keys,
28                                wordset,
29                                word_sequences,
30                                tagset,
31                                tag_sequences)
32
33     def __len__(self):
34         return len(self.sentences)
35
36     def __iter__(self):
37         return iter(self.sentences.items())
38
39 ''' 7. SPLIT DATA INTO TRAINING & TEST SETS Training : Test = 0.8:0.2'''
40 key_list = list(keys)
41 train_ratio=0.8
42 random.shuffle(key_list)
43 split = int(train_ratio * len(key_list))
44 training_data = Subset(sentences, key_list[:split])
45 test_data = Subset(sentences, key_list[split:])
46
47 print("Training set has {} sentences.".format(len(training_data.keys)))
48 print("Test set has {} sentences.\n".format(len(test_data.keys)))
49
50

```

Training set has 45872 sentences.

Test set has 11468 sentences.

In [4]:

```

1 #Data Description
2
3 ##### 2. COUNT NUMBER OF EACH TAG IN ENTIRE CORPUS #####
4 ##### single_tag_counts[tag] = k #####
5 all_sentences = [list(sentence) for sentence in training_data.Y]
6 all_tags = chain.from_iterable(all_sentences)
7 single_tag_counts = dict(Counter(all_tags))
8
9 ##### 3. COUNT NUMBER OF EACH TAG PAIR IN ENTIRE CORPUS #####
10 ##### pair_tag_counts[(tag_1, tag_2)] = k #####
11 sentences = [s for s in all_sentences if len(s) > 1] # discard any sequences o
12 pairs = []
13 for s in sentences:
14     pairs.extend([(s[i-1], s[i]) for i in range(1, len(s))])
15
16 pair_tag_counts = dict(Counter(pairs))
17
18 if len(pair_tag_counts) < len(training_data.tagset)**2:
19     for tag1 in training_data.tagset:
20         for tag2 in training_data.tagset:
21             if (tag1, tag2) not in pair_tag_counts:
22                 pair_tag_counts[(tag1, tag2)] = 0
23
24 ## 4. COUNT NUMBER OF EACH TAG APPEARING IN THE BEGINNING or END OF SENTENCE ##
25 ##### start_tag_counts[tag] = k #####
26 ##### end_tag_counts[tag] = k #####
27
28 start_tag_counts = dict(Counter([sentence[0] for sentence in training_data.Y]))
29 end_tag_counts = dict(Counter([sentence[-1] for sentence in training_data.Y]))
30
31 ### if any tag has NO sentences starting/ending with it, set its value to 0:
32 if len(start_tag_counts) < len(training_data.tagset):
33     for tag in training_data.tagset:
34         if tag not in start_tag_counts:
35             start_tag_counts[tag] = 0
36
37 if len(end_tag_counts) < len(training_data.tagset):
38     for tag in training_data.tagset:
39         if tag not in end_tag_counts:
40             end_tag_counts[tag] = 0
41
42 ##### 5. COUNT NUMBER OF (TAG_i, WORD_i) PAIRS #####
43 ##### pair_counts[tag][word] = k #####
44 pair_counts = defaultdict(lambda: defaultdict(lambda : 0))
45
46 for sentence_idx, sentence in enumerate(training_data.Y):
47     for word_idx, tag in enumerate(sentence):
48         word = training_data.X[sentence_idx][word_idx]
49         pair_counts[tag][word] += 1
50

```

In [5]:

```

1 #HMM Model Goes Here
2
3 ##### 6. BUILD HMM MODEL #####
4 HMM_model = HiddenMarkovModel(name = "HMM-Tagger")
5 tag_states = [] # state for each tag
6
7 ##### (6.1) ADD STATES w/ EMISSION PROBABILITIES #####
8
9 for tag in training_data.tagset:
10     tag_emissions = DiscreteDistribution({word:pair_counts[tag][word]/single_tag_counts[tag]
11                                           for word in training_data.vocab})
12     tag_state = State(tag_emissions, name = tag)
13     tag_states.append(tag_state)
14     HMM_model.add_states(tag_state)
15
16 ##### (6.2) ADD TRANSITIONS w/ TRANSITION PROBABILITIES #####
17
18 n_sentences = len(training_data.keys)
19
20 for tag_state1 in tag_states:
21     for tag_state2 in tag_states:
22         tag1, tag2 = tag_state1.name, tag_state2.name
23         HMM_model.add_transition(HMM_model.start, tag_state1, start_tag_counts[tag1])
24         HMM_model.add_transition(tag_state1, HMM_model.end, end_tag_counts[tag1])
25         HMM_model.add_transition(tag_state1, tag_state2, pair_tag_counts[(tag1, tag2)])
26
27 HMM_model.bake()

```

In [6]:

```

1 #Model Accuracy Evaluation
2
3 ##### 7. MAKE PREDICTIONS ON TRAINING SET #####
4 ##### DISPLAY THE OUTPUT ON CONSOLE #####
5
6 train_correct = 0 # number of correct predictions so far
7 train_count = 0   # number of predictions so far
8 print_i = 100
9
10 ### ITERATE PER SENTENCE
11 for words, true_tags in zip(training_data.X, training_data.Y):
12     try:
13         # Viterbi Path: most likely sequence of STATES that generated the sequence
14         _, viterbi_path = HMM_model.viterbi([w for w in words])
15         predicted_tags = [state[1].name for state in viterbi_path[1:-1]]
16         train_correct += sum(pred == true for pred, true in zip(predicted_tags,
17
18             if print_i == 100: # print a sample result
19                 print("Training Sentence: \n", words)
20                 print()
21                 print("Predicted Tags: \n", predicted_tags)
22                 print()
23                 print("True Tags: \n", true_tags)
24                 print_i += 1
25     except:
26         pass
27     train_count += len(words)
28
29 train_acc = train_correct/train_count
30 print("\nTraining Accuracy: {:.2f}%".format(100 * train_acc))
31 print()
32 print()

```

Training Sentence:

```
('And', 'Early', 'Spring', 'seized', 'the', 'whip', 'and', 'said',
':')
```

Predicted Tags:

```
['CONJ', 'ADJ', 'NOUN', 'VERB', 'DET', 'NOUN', 'CONJ', 'VERB', '.']
```

True Tags:

```
('CONJ', 'ADJ', 'NOUN', 'VERB', 'DET', 'NOUN', 'CONJ', 'VERB', '.')
```

Training Accuracy: 97.55%

In [110]:

```

1 #Adds code blocks wherever you feel necessary
2
3 ##### 8. MAKE PREDICTIONS ON TEST SET #####
4 ##### DISPLAY THE OUTPUT ON CONSOLE #####
5 test_correct = 0
6 test_count = 0
7
8 ### ITERATE PER SENTENCE
9 for words, true_tags in zip(test_data.X, test_data.Y):
10     try:
11         # Only consider words contained in training set's vocab
12         _, viterbi_path = HMM_model.viterbi([w if w in training_data.vocab else
13         predicted_tags = [state[1].name for state in viterbi_path[1:-1]]
14         test_correct += sum(pred == true for pred, true in zip(predicted_tags, true_tags))
15
16         if print_i == 101: # print a sample result
17             print("Test Sentence: \n", words)
18             print()
19             print("Predicted Tags: \n", predicted_tags)
20             print()
21             print("True Tags: \n", true_tags)
22             print_i += 1
23     except:
24         pass
25     test_count += len(words)
26
27 test_acc = test_correct/test_count
28 print("\nTest Accuracy: {:.2f}%".format(100 * test_acc))

```

Test Sentence:

```

('Police', 'said', 'he', 'became', 'ill', 'while', 'parked', 'in', 'front', 'of', 'a', 'barber', 'shop', 'at', '229', 'West', 'Pratt', 'Street', '.')

```

Predicted Tags:

```

['NOUN', 'VERB', 'PRON', 'VERB', 'ADJ', 'NOUN', 'VERB', 'ADP', 'NOUN', 'ADP', 'DET', 'NOUN', 'NOUN', 'ADP', 'NUM', 'ADJ', 'NOUN', 'NOUN', '.']

```

True Tags:

```

('NOUN', 'VERB', 'PRON', 'VERB', 'ADJ', 'ADP', 'VERB', 'ADP', 'NOUN', 'ADP', 'DET', 'NOUN', 'NOUN', 'ADP', 'NUM', 'ADJ', 'NOUN', 'NOUN', '.')

```

Test Accuracy: 95.98%

Happy Coding!