

# app

October 8, 2019

## POS Tagging with HMM and Sentence Generation

The training dataset is a subset of the Brown corpus, where each file contains sentences in the form of tokenized words followed by POS tags. Each line contains one sentence. Training dataset can be downloaded from here: <https://bit.ly/2kJI0yc> The test dataset (which is another subset of the Brown corpus, containing tokenized words but no tags) can be downloaded from here: <https://bit.ly/2lMybzP> Information regarding the categories of the dataset can be found at: <https://bit.ly/2mhF6RT>.

Your task is to implement a part-of-speech tagger using a bi-gram HMM. Given an observation sequence of  $n$  words  $w_{n1}$ , choose the most probable sequence of POS tags  $t_{n1}$ . For the questions below, please submit both code and output.

[Note: During training, for a word to be counted as unknown, the frequency of the word in training set should not exceed a threshold (e.g. 5). You can pick a threshold based on your algorithm design. Also, you can implement smoothing technique based on your own choice, e.g. add-.]

```
In [1]: import glob
import re
import math
import random
import pandas as pd
from collections import Counter
from nltk.util import ngrams
from nltk.tokenize import sent_tokenize, word_tokenize

In [2]: def load_data(input_directory):
    word_tokens = []
    tag_tokens = []
    word_tag_tokens = []
    transition_tag_tokens = []
    sentence_count = 0
    for file_name in glob.glob(input_directory + "/*"):
        print("Preprocessing: {}".format(file_name))
        file_pointer = open(file_name, "r")

        for line in file_pointer:

            # Remove duplicate spaces
```

```

file_line_content = re.sub(' +', ' ', line)

# Remove new line characters
file_line_content = line.replace("\n", " ")

# Strip of begin and end spaces
file_line_content = file_line_content.strip()

# If line is not empty
if file_line_content != "":
    sentence_count = sentence_count + 1
    line_content_list = file_line_content.split(" ")

    # Append start tag
    transition_tag_tokens.append('START')
    print(line_content_list)
    for i in line_content_list:
        word_tag_tokens.append(i)
        split_tokens = i.split('/')
        word = split_tokens[0]
        tag = split_tokens[-1]
        word_tokens.append(word)
        tag_tokens.append(tag)
        transition_tag_tokens.append(tag)

    # Append end tag
    transition_tag_tokens.append('END')

file_pointer.close()

return sentence_count, word_tokens, tag_tokens, word_tag_tokens, transition_tag_tokens

In [3]: def replace_low_count_words(word_tokens, cut_off_count, word_tag_tokens):
    word_tokens_with_count = Counter(word_tokens)
    candidate_words = {}
    for word in word_tokens_with_count:
        if word_tokens_with_count[word] <= cut_off_count:
            candidate_words[word] = 1

    ## Remove all the words <= cut_off_count
    for i in range(len(word_tokens)):
        if word_tokens[i] in candidate_words:
            word_tokens[i] = 'UNK'
            split_tokens = word_tag_tokens[i].split('/')
            word = split_tokens[0]
            tag = split_tokens[-1]
            word_tag_tokens[i] = 'UNK' + '/' + tag

```

```

        return word_tokens, word_tag_tokens

In [4]: print("#### Loading train date")
        # train_sentence_count, train_word_tokens, train_tag_tokens, train_word_tag_tokens, \
        #     train_transition_tag_tokens = load_data('../input/custom_train')

        train_sentence_count, train_word_tokens, train_tag_tokens, train_word_tag_tokens, \
        train_transition_tag_tokens = load_data('../input/brown_train')

        print(" Number of Sentences: {}, Word list count: {}, Tag list count: {}, Transition T
              " Word Tag List"
              .format(train_sentence_count, len(train_word_tokens), len(train_tag_tokens), len
              )

        print("### Replace word tokens <= 5 with 'UNK'")
        train_word_tokens, train_word_tag_tokens = replace_low_count_words(train_word_tokens,
        print("Number of tokens after replacement: Word - {}, Word_Tag - {}".format(len(train_

#### Loading train date
Preprocessing: ../input/brown_train/cd05
Preprocessing: ../input/brown_train/cf37
Preprocessing: ../input/brown_train/cf08
Preprocessing: ../input/brown_train/cl06
Preprocessing: ../input/brown_train/cj68
Preprocessing: ../input/brown_train/cf30
Preprocessing: ../input/brown_train/cd02
Preprocessing: ../input/brown_train/cl01
Preprocessing: ../input/brown_train/cj57
Preprocessing: ../input/brown_train/cf06
Preprocessing: ../input/brown_train/cl08
Preprocessing: ../input/brown_train/cj61
Preprocessing: ../input/brown_train/cf39
Preprocessing: ../input/brown_train/cn05
Preprocessing: ../input/brown_train/cj59
Preprocessing: ../input/brown_train/cf01
Preprocessing: ../input/brown_train/cn02
Preprocessing: ../input/brown_train/cj66
Preprocessing: ../input/brown_train/cj32
Preprocessing: ../input/brown_train/ch07
Preprocessing: ../input/brown_train/cj35
Preprocessing: ../input/brown_train/cb09
Preprocessing: ../input/brown_train/cj03
Preprocessing: ../input/brown_train/cj04
Preprocessing: ../input/brown_train/ch09
Preprocessing: ../input/brown_train/cb07
Preprocessing: ../input/brown_train/cj58
Preprocessing: ../input/brown_train/cn03
Preprocessing: ../input/brown_train/cj67

```

Preprocessing: ../input/brown\_train/cl09  
Preprocessing: ../input/brown\_train/cf07  
Preprocessing: ../input/brown\_train/cj60  
Preprocessing: ../input/brown\_train/cn04  
Preprocessing: ../input/brown\_train/cf38  
Preprocessing: ../input/brown\_train/cj69  
Preprocessing: ../input/brown\_train/cd03  
Preprocessing: ../input/brown\_train/cf31  
Preprocessing: ../input/brown\_train/cj56  
Preprocessing: ../input/brown\_train/cf36  
Preprocessing: ../input/brown\_train/cd04  
Preprocessing: ../input/brown\_train/cj51  
Preprocessing: ../input/brown\_train/cl07  
Preprocessing: ../input/brown\_train/cf09  
Preprocessing: ../input/brown\_train/cj05  
Preprocessing: ../input/brown\_train/cb06  
Preprocessing: ../input/brown\_train/cj02  
Preprocessing: ../input/brown\_train/ch30  
Preprocessing: ../input/brown\_train/cb01  
Preprocessing: ../input/brown\_train/cb08  
Preprocessing: ../input/brown\_train/cj34  
Preprocessing: ../input/brown\_train/ch06  
Preprocessing: ../input/brown\_train/ch01  
Preprocessing: ../input/brown\_train/cj33  
Preprocessing: ../input/brown\_train/cr03  
Preprocessing: ../input/brown\_train/ce14  
Preprocessing: ../input/brown\_train/cg26  
Preprocessing: ../input/brown\_train/cg19  
Preprocessing: ../input/brown\_train/cg21  
Preprocessing: ../input/brown\_train/ce13  
Preprocessing: ../input/brown\_train/cg17  
Preprocessing: ../input/brown\_train/ce25  
Preprocessing: ../input/brown\_train/ca41  
Preprocessing: ../input/brown\_train/cg28  
Preprocessing: ../input/brown\_train/cp07  
Preprocessing: ../input/brown\_train/ce22  
Preprocessing: ../input/brown\_train/cg10  
Preprocessing: ../input/brown\_train/ck23  
Preprocessing: ../input/brown\_train/ca12  
Preprocessing: ../input/brown\_train/cg44  
Preprocessing: ../input/brown\_train/ck24  
Preprocessing: ../input/brown\_train/cg43  
Preprocessing: ../input/brown\_train/ca15  
Preprocessing: ../input/brown\_train/ck12  
Preprocessing: ../input/brown\_train/cc11  
Preprocessing: ../input/brown\_train/ca23  
Preprocessing: ../input/brown\_train/cg75  
Preprocessing: ../input/brown\_train/cg72

Preprocessing: ../input/brown\_train/ca24  
Preprocessing: ../input/brown\_train/cc16  
Preprocessing: ../input/brown\_train/cg11  
Preprocessing: ../input/brown\_train/ce23  
Preprocessing: ../input/brown\_train/cp06  
Preprocessing: ../input/brown\_train/cp01  
Preprocessing: ../input/brown\_train/ce24  
Preprocessing: ../input/brown\_train/cg16  
Preprocessing: ../input/brown\_train/ca40  
Preprocessing: ../input/brown\_train/cg29  
Preprocessing: ../input/brown\_train/cr05  
Preprocessing: ../input/brown\_train/ce12  
Preprocessing: ../input/brown\_train/cg20  
Preprocessing: ../input/brown\_train/cp08  
Preprocessing: ../input/brown\_train/cg27  
Preprocessing: ../input/brown\_train/ce15  
Preprocessing: ../input/brown\_train/cr02  
Preprocessing: ../input/brown\_train/cg18  
Preprocessing: ../input/brown\_train/ck14  
Preprocessing: ../input/brown\_train/cg73  
Preprocessing: ../input/brown\_train/cc17  
Preprocessing: ../input/brown\_train/ca25  
Preprocessing: ../input/brown\_train/ck13  
Preprocessing: ../input/brown\_train/ca22  
Preprocessing: ../input/brown\_train/cc10  
Preprocessing: ../input/brown\_train/cg74  
Preprocessing: ../input/brown\_train/ck25  
Preprocessing: ../input/brown\_train/cg42  
Preprocessing: ../input/brown\_train/ca14  
Preprocessing: ../input/brown\_train/ck22  
Preprocessing: ../input/brown\_train/ca13  
Preprocessing: ../input/brown\_train/cg45  
Preprocessing: ../input/brown\_train/cg67  
Preprocessing: ../input/brown\_train/ca31  
Preprocessing: ../input/brown\_train/cc03  
Preprocessing: ../input/brown\_train/cg58  
Preprocessing: ../input/brown\_train/cc04  
Preprocessing: ../input/brown\_train/ca36  
Preprocessing: ../input/brown\_train/cg60  
Preprocessing: ../input/brown\_train/ca09  
Preprocessing: ../input/brown\_train/ck07  
Preprocessing: ../input/brown\_train/cg56  
Preprocessing: ../input/brown\_train/cg69  
Preprocessing: ../input/brown\_train/ca07  
Preprocessing: ../input/brown\_train/ck09  
Preprocessing: ../input/brown\_train/cg51  
Preprocessing: ../input/brown\_train/ca38  
Preprocessing: ../input/brown\_train/ce08

Preprocessing: ../input/brown\_train/cg05  
Preprocessing: ../input/brown\_train/cp12  
Preprocessing: ../input/brown\_train/cp15  
Preprocessing: ../input/brown\_train/cg02  
Preprocessing: ../input/brown\_train/ce30  
Preprocessing: ../input/brown\_train/cm05  
Preprocessing: ../input/brown\_train/cp23  
Preprocessing: ../input/brown\_train/cg34  
Preprocessing: ../input/brown\_train/ce06  
Preprocessing: ../input/brown\_train/cm02  
Preprocessing: ../input/brown\_train/ce01  
Preprocessing: ../input/brown\_train/cg33  
Preprocessing: ../input/brown\_train/cp24  
Preprocessing: ../input/brown\_train/ck08  
Preprocessing: ../input/brown\_train/ca06  
Preprocessing: ../input/brown\_train/cg50  
Preprocessing: ../input/brown\_train/ca39  
Preprocessing: ../input/brown\_train/ca01  
Preprocessing: ../input/brown\_train/cg68  
Preprocessing: ../input/brown\_train/ca37  
Preprocessing: ../input/brown\_train/cc05  
Preprocessing: ../input/brown\_train/cg61  
Preprocessing: ../input/brown\_train/ck06  
Preprocessing: ../input/brown\_train/ca08  
Preprocessing: ../input/brown\_train/cg66  
Preprocessing: ../input/brown\_train/cc02  
Preprocessing: ../input/brown\_train/ca30  
Preprocessing: ../input/brown\_train/ck01  
Preprocessing: ../input/brown\_train/cg59  
Preprocessing: ../input/brown\_train/cp25  
Preprocessing: ../input/brown\_train/cg32  
Preprocessing: ../input/brown\_train/cm04  
Preprocessing: ../input/brown\_train/ce07  
Preprocessing: ../input/brown\_train/cg35  
Preprocessing: ../input/brown\_train/cp22  
Preprocessing: ../input/brown\_train/ce31  
Preprocessing: ../input/brown\_train/cg03  
Preprocessing: ../input/brown\_train/cp14  
Preprocessing: ../input/brown\_train/ce09  
Preprocessing: ../input/brown\_train/cp13  
Preprocessing: ../input/brown\_train/cg04  
Preprocessing: ../input/brown\_train/ce36  
Preprocessing: ../input/brown\_train/cb12  
Preprocessing: ../input/brown\_train/cj11  
Preprocessing: ../input/brown\_train/ch23  
Preprocessing: ../input/brown\_train/cb15  
Preprocessing: ../input/brown\_train/cj29  
Preprocessing: ../input/brown\_train/ch24

Preprocessing: ../input/brown\_train/cj16  
Preprocessing: ../input/brown\_train/cf47  
Preprocessing: ../input/brown\_train/cb23  
Preprocessing: ../input/brown\_train/ch12  
Preprocessing: ../input/brown\_train/cj20  
Preprocessing: ../input/brown\_train/cb24  
Preprocessing: ../input/brown\_train/cj18  
Preprocessing: ../input/brown\_train/cf40  
Preprocessing: ../input/brown\_train/cj27  
Preprocessing: ../input/brown\_train/ch15  
Preprocessing: ../input/brown\_train/cj73  
Preprocessing: ../input/brown\_train/cn28  
Preprocessing: ../input/brown\_train/cj80  
Preprocessing: ../input/brown\_train/cj74  
Preprocessing: ../input/brown\_train/cl22  
Preprocessing: ../input/brown\_train/cn10  
Preprocessing: ../input/brown\_train/cf13  
Preprocessing: ../input/brown\_train/cl14  
Preprocessing: ../input/brown\_train/cn26  
Preprocessing: ../input/brown\_train/cj42  
Preprocessing: ../input/brown\_train/cn19  
Preprocessing: ../input/brown\_train/cf25  
Preprocessing: ../input/brown\_train/cd17  
Preprocessing: ../input/brown\_train/cj45  
Preprocessing: ../input/brown\_train/cn21  
Preprocessing: ../input/brown\_train/cl13  
Preprocessing: ../input/brown\_train/cd10  
Preprocessing: ../input/brown\_train/cf22  
Preprocessing: ../input/brown\_train/cj19  
Preprocessing: ../input/brown\_train/cf41  
Preprocessing: ../input/brown\_train/ch14  
Preprocessing: ../input/brown\_train/cj26  
Preprocessing: ../input/brown\_train/cf46  
Preprocessing: ../input/brown\_train/cb22  
Preprocessing: ../input/brown\_train/cj21  
Preprocessing: ../input/brown\_train/ch13  
Preprocessing: ../input/brown\_train/cj28  
Preprocessing: ../input/brown\_train/cb14  
Preprocessing: ../input/brown\_train/cj17  
Preprocessing: ../input/brown\_train/ch25  
Preprocessing: ../input/brown\_train/cb13  
Preprocessing: ../input/brown\_train/ch22  
Preprocessing: ../input/brown\_train/cj10  
Preprocessing: ../input/brown\_train/cf48  
Preprocessing: ../input/brown\_train/cj44  
Preprocessing: ../input/brown\_train/cl12  
Preprocessing: ../input/brown\_train/cn20  
Preprocessing: ../input/brown\_train/cf23

Preprocessing: ../input/brown\_train/cd11  
Preprocessing: ../input/brown\_train/cn27  
Preprocessing: ../input/brown\_train/cl15  
Preprocessing: ../input/brown\_train/cj43  
Preprocessing: ../input/brown\_train/cd16  
Preprocessing: ../input/brown\_train/cf24  
Preprocessing: ../input/brown\_train/cn18  
Preprocessing: ../input/brown\_train/cj75  
Preprocessing: ../input/brown\_train/cn11  
Preprocessing: ../input/brown\_train/cl23  
Preprocessing: ../input/brown\_train/cf12  
Preprocessing: ../input/brown\_train/cl24  
Preprocessing: ../input/brown\_train/cn16  
Preprocessing: ../input/brown\_train/cj72  
Preprocessing: ../input/brown\_train/cf15  
Preprocessing: ../input/brown\_train/cn29  
Preprocessing: ../input/brown\_train/cj36  
Preprocessing: ../input/brown\_train/ch04  
Preprocessing: ../input/brown\_train/cj09  
Preprocessing: ../input/brown\_train/ch03  
Preprocessing: ../input/brown\_train/cj31  
Preprocessing: ../input/brown\_train/cj07  
Preprocessing: ../input/brown\_train/cj38  
Preprocessing: ../input/brown\_train/cb04  
Preprocessing: ../input/brown\_train/cb03  
Preprocessing: ../input/brown\_train/cd01  
Preprocessing: ../input/brown\_train/cf33  
Preprocessing: ../input/brown\_train/cl02  
Preprocessing: ../input/brown\_train/cj54  
Preprocessing: ../input/brown\_train/cf34  
Preprocessing: ../input/brown\_train/cn08  
Preprocessing: ../input/brown\_train/cd06  
Preprocessing: ../input/brown\_train/cj53  
Preprocessing: ../input/brown\_train/cl05  
Preprocessing: ../input/brown\_train/cf02  
Preprocessing: ../input/brown\_train/cn01  
Preprocessing: ../input/brown\_train/cj65  
Preprocessing: ../input/brown\_train/cf05  
Preprocessing: ../input/brown\_train/cj62  
Preprocessing: ../input/brown\_train/cn06  
Preprocessing: ../input/brown\_train/cd08  
Preprocessing: ../input/brown\_train/cj01  
Preprocessing: ../input/brown\_train/cb02  
Preprocessing: ../input/brown\_train/cj06  
Preprocessing: ../input/brown\_train/cb05  
Preprocessing: ../input/brown\_train/cj39  
Preprocessing: ../input/brown\_train/cj30  
Preprocessing: ../input/brown\_train/ch02



Preprocessing: ../input/brown\_train/ch05  
Preprocessing: ../input/brown\_train/cj37  
Preprocessing: ../input/brown\_train/cj08  
Preprocessing: ../input/brown\_train/cf04  
Preprocessing: ../input/brown\_train/cj63  
Preprocessing: ../input/brown\_train/cd09  
Preprocessing: ../input/brown\_train/cn07  
Preprocessing: ../input/brown\_train/cf03  
Preprocessing: ../input/brown\_train/cj64  
Preprocessing: ../input/brown\_train/cd07  
Preprocessing: ../input/brown\_train/cn09  
Preprocessing: ../input/brown\_train/cf35  
Preprocessing: ../input/brown\_train/cj52  
Preprocessing: ../input/brown\_train/cl04  
Preprocessing: ../input/brown\_train/cf32  
Preprocessing: ../input/brown\_train/cl03  
Preprocessing: ../input/brown\_train/cj55  
Preprocessing: ../input/brown\_train/ck27  
Preprocessing: ../input/brown\_train/ca29  
Preprocessing: ../input/brown\_train/cg40  
Preprocessing: ../input/brown\_train/ca16  
Preprocessing: ../input/brown\_train/ck18  
Preprocessing: ../input/brown\_train/ck20  
Preprocessing: ../input/brown\_train/ca11  
Preprocessing: ../input/brown\_train/cg47  
Preprocessing: ../input/brown\_train/ca18  
Preprocessing: ../input/brown\_train/ck16  
Preprocessing: ../input/brown\_train/cg71  
Preprocessing: ../input/brown\_train/cc15  
Preprocessing: ../input/brown\_train/ck29  
Preprocessing: ../input/brown\_train/ca27  
Preprocessing: ../input/brown\_train/cg49  
Preprocessing: ../input/brown\_train/ck11  
Preprocessing: ../input/brown\_train/ca20  
Preprocessing: ../input/brown\_train/cc12  
Preprocessing: ../input/brown\_train/ce10  
Preprocessing: ../input/brown\_train/cg22  
Preprocessing: ../input/brown\_train/cr07  
Preprocessing: ../input/brown\_train/cg25  
Preprocessing: ../input/brown\_train/ce17  
Preprocessing: ../input/brown\_train/ce28  
Preprocessing: ../input/brown\_train/cp04  
Preprocessing: ../input/brown\_train/cg13  
Preprocessing: ../input/brown\_train/ce21  
Preprocessing: ../input/brown\_train/cr09  
Preprocessing: ../input/brown\_train/ce26  
Preprocessing: ../input/brown\_train/cg14  
Preprocessing: ../input/brown\_train/ca42

Preprocessing: ../input/brown\_train/cp03  
Preprocessing: ../input/brown\_train/cg48  
Preprocessing: ../input/brown\_train/ck10  
Preprocessing: ../input/brown\_train/cc13  
Preprocessing: ../input/brown\_train/ca21  
Preprocessing: ../input/brown\_train/ck17  
Preprocessing: ../input/brown\_train/ca19  
Preprocessing: ../input/brown\_train/cg70  
Preprocessing: ../input/brown\_train/ca26  
Preprocessing: ../input/brown\_train/ck28  
Preprocessing: ../input/brown\_train/cc14  
Preprocessing: ../input/brown\_train/ck21  
Preprocessing: ../input/brown\_train/ca10  
Preprocessing: ../input/brown\_train/cg46  
Preprocessing: ../input/brown\_train/ca28  
Preprocessing: ../input/brown\_train/ck26  
Preprocessing: ../input/brown\_train/cg41  
Preprocessing: ../input/brown\_train/ck19  
Preprocessing: ../input/brown\_train/ca17  
Preprocessing: ../input/brown\_train/cp02  
Preprocessing: ../input/brown\_train/cg15  
Preprocessing: ../input/brown\_train/ce27  
Preprocessing: ../input/brown\_train/ca43  
Preprocessing: ../input/brown\_train/ce18  
Preprocessing: ../input/brown\_train/ca44  
Preprocessing: ../input/brown\_train/ce20  
Preprocessing: ../input/brown\_train/cg12  
Preprocessing: ../input/brown\_train/cp05  
Preprocessing: ../input/brown\_train/cr08  
Preprocessing: ../input/brown\_train/ce16  
Preprocessing: ../input/brown\_train/cg24  
Preprocessing: ../input/brown\_train/cr01  
Preprocessing: ../input/brown\_train/ce29  
Preprocessing: ../input/brown\_train/cr06  
Preprocessing: ../input/brown\_train/cg23  
Preprocessing: ../input/brown\_train/ce11  
Preprocessing: ../input/brown\_train/cp29  
Preprocessing: ../input/brown\_train/cp16  
Preprocessing: ../input/brown\_train/ce33  
Preprocessing: ../input/brown\_train/cg01  
Preprocessing: ../input/brown\_train/cg39  
Preprocessing: ../input/brown\_train/cg06  
Preprocessing: ../input/brown\_train/ce34  
Preprocessing: ../input/brown\_train/cp11  
Preprocessing: ../input/brown\_train/cm01  
Preprocessing: ../input/brown\_train/cp18  
Preprocessing: ../input/brown\_train/cg30  
Preprocessing: ../input/brown\_train/ce02

Preprocessing: ../input/brown\_train/cp27  
Preprocessing: ../input/brown\_train/cg08  
Preprocessing: ../input/brown\_train/cm06  
Preprocessing: ../input/brown\_train/cp20  
Preprocessing: ../input/brown\_train/ce05  
Preprocessing: ../input/brown\_train/cg37  
Preprocessing: ../input/brown\_train/ca35  
Preprocessing: ../input/brown\_train/cg63  
Preprocessing: ../input/brown\_train/ck04  
Preprocessing: ../input/brown\_train/cg64  
Preprocessing: ../input/brown\_train/ca32  
Preprocessing: ../input/brown\_train/ck03  
Preprocessing: ../input/brown\_train/ca04  
Preprocessing: ../input/brown\_train/cg52  
Preprocessing: ../input/brown\_train/cc09  
Preprocessing: ../input/brown\_train/cg55  
Preprocessing: ../input/brown\_train/ca03  
Preprocessing: ../input/brown\_train/cg09  
Preprocessing: ../input/brown\_train/cg36  
Preprocessing: ../input/brown\_train/ce04  
Preprocessing: ../input/brown\_train/cp21  
Preprocessing: ../input/brown\_train/cp19  
Preprocessing: ../input/brown\_train/cp26  
Preprocessing: ../input/brown\_train/ce03  
Preprocessing: ../input/brown\_train/cg31  
Preprocessing: ../input/brown\_train/cg38  
Preprocessing: ../input/brown\_train/cp10  
Preprocessing: ../input/brown\_train/ce35  
Preprocessing: ../input/brown\_train/cg07  
Preprocessing: ../input/brown\_train/cp28  
Preprocessing: ../input/brown\_train/ce32  
Preprocessing: ../input/brown\_train/cp17  
Preprocessing: ../input/brown\_train/cg54  
Preprocessing: ../input/brown\_train/ca02  
Preprocessing: ../input/brown\_train/ca05  
Preprocessing: ../input/brown\_train/cg53  
Preprocessing: ../input/brown\_train/cc08  
Preprocessing: ../input/brown\_train/cg65  
Preprocessing: ../input/brown\_train/ca33  
Preprocessing: ../input/brown\_train/cc01  
Preprocessing: ../input/brown\_train/ck02  
Preprocessing: ../input/brown\_train/cc06  
Preprocessing: ../input/brown\_train/cg62  
Preprocessing: ../input/brown\_train/ck05  
Preprocessing: ../input/brown\_train/cj77  
Preprocessing: ../input/brown\_train/cn13  
Preprocessing: ../input/brown\_train/cf10  
Preprocessing: ../input/brown\_train/cj48

Preprocessing: ../input/brown\_train/cn14  
Preprocessing: ../input/brown\_train/cf28  
Preprocessing: ../input/brown\_train/cj70  
Preprocessing: ../input/brown\_train/cf17  
Preprocessing: ../input/brown\_train/cl19  
Preprocessing: ../input/brown\_train/cj46  
Preprocessing: ../input/brown\_train/cl10  
Preprocessing: ../input/brown\_train/cn22  
Preprocessing: ../input/brown\_train/cf21  
Preprocessing: ../input/brown\_train/cd13  
Preprocessing: ../input/brown\_train/cj79  
Preprocessing: ../input/brown\_train/cn25  
Preprocessing: ../input/brown\_train/cf19  
Preprocessing: ../input/brown\_train/cl17  
Preprocessing: ../input/brown\_train/cj41  
Preprocessing: ../input/brown\_train/cf26  
Preprocessing: ../input/brown\_train/ch18  
Preprocessing: ../input/brown\_train/cb16  
Preprocessing: ../input/brown\_train/cj15  
Preprocessing: ../input/brown\_train/ch27  
Preprocessing: ../input/brown\_train/cb11  
Preprocessing: ../input/brown\_train/ch20  
Preprocessing: ../input/brown\_train/cj12  
Preprocessing: ../input/brown\_train/cb27  
Preprocessing: ../input/brown\_train/ch29  
Preprocessing: ../input/brown\_train/cf43  
Preprocessing: ../input/brown\_train/ch16  
Preprocessing: ../input/brown\_train/cb18  
Preprocessing: ../input/brown\_train/cj24  
Preprocessing: ../input/brown\_train/cf44  
Preprocessing: ../input/brown\_train/cb20  
Preprocessing: ../input/brown\_train/cj23  
Preprocessing: ../input/brown\_train/ch11  
Preprocessing: ../input/brown\_train/cl16  
Preprocessing: ../input/brown\_train/cf18  
Preprocessing: ../input/brown\_train/cn24  
Preprocessing: ../input/brown\_train/cj40  
Preprocessing: ../input/brown\_train/cf27  
Preprocessing: ../input/brown\_train/cd15  
Preprocessing: ../input/brown\_train/cj47  
Preprocessing: ../input/brown\_train/cn23  
Preprocessing: ../input/brown\_train/cl11  
Preprocessing: ../input/brown\_train/cd12  
Preprocessing: ../input/brown\_train/cf20  
Preprocessing: ../input/brown\_train/cj78  
Preprocessing: ../input/brown\_train/cf29  
Preprocessing: ../input/brown\_train/cn15  
Preprocessing: ../input/brown\_train/cj71

```

Preprocessing: ../input/brown_train/cl18
Preprocessing: ../input/brown_train/cf16
Preprocessing: ../input/brown_train/cj76
Preprocessing: ../input/brown_train/cl20
Preprocessing: ../input/brown_train/cn12
Preprocessing: ../input/brown_train/cf11
Preprocessing: ../input/brown_train/cj49
Preprocessing: ../input/brown_train/cf45
Preprocessing: ../input/brown_train/cb21
Preprocessing: ../input/brown_train/ch10
Preprocessing: ../input/brown_train/cj22
Preprocessing: ../input/brown_train/ch28
Preprocessing: ../input/brown_train/cb26
Preprocessing: ../input/brown_train/cf42
Preprocessing: ../input/brown_train/cj25
Preprocessing: ../input/brown_train/cb19
Preprocessing: ../input/brown_train/ch17
Preprocessing: ../input/brown_train/cb10
Preprocessing: ../input/brown_train/cj13
Preprocessing: ../input/brown_train/ch21
Preprocessing: ../input/brown_train/cb17
Preprocessing: ../input/brown_train/ch19
Preprocessing: ../input/brown_train/ch26
Preprocessing: ../input/brown_train/cj14

```

```

Number of Sentences: 55684, Word list count: 1126281, Tag list count: 1126281, Transition Tag
### Replace word tokens <= 5 with 'UNK'
Number of tokens after replacement: Word - 1126281, Word_Tag - 1126281

```

4.1 Obtain frequency counts from the collection of all the training files (counted together). You will need the following types of frequency counts: word-tag counts, tag un-igram counts, and tag bigram counts. Let's denote these by  $C(w_i, t_i)$ ,  $C(t_i)$  and  $C(t_{i1}, t_i)$  respectively. Report these quantities in different output files.

```

In [5]: def get_unigrams(token_list):
        tokens_with_count = Counter(token_list)
        return tokens_with_count

In [6]: def get_ngrams(token_list, n):
        ngrams_zip = zip(*[token_list[i:] for i in range(n)])
        ngrams_list = [" ".join(element) for element in ngrams_zip]
        ngrams_keys_counts = Counter(ngrams_list)
        return ngrams_keys_counts

In [7]: print("Unigrams for words")
        train_word_unigrams = get_unigrams(train_word_tokens)
        print(len(train_word_unigrams))
        vocab_size = len(train_word_unigrams)
        print("Vocab size {}".format(vocab_size))

```

```

print("Unigrams for word-tag")
print(len(train_word_tag_tokens))
train_word_tag_unigrams = get_unigrams(train_word_tag_tokens)
with open('../output/word_tag_unigrams.txt', 'w') as word_tag_unigrams_output_file:
    word_tag_unigrams_output_file.write(str(train_word_tag_unigrams))
print(len(train_word_tag_unigrams))

print("Unigrams for tags")
train_tag_unigrams = get_unigrams(train_tag_tokens)
print(len(train_tag_unigrams))

print("Unigrams for transition tags")
train_transition_tag_unigrams = get_unigrams(train_transition_tag_tokens)
with open('../output/tag_unigrams.txt', 'w') as transition_tag_unigrams_output_file:
    transition_tag_unigrams_output_file.write(str(train_transition_tag_unigrams))
print(len(train_transition_tag_unigrams))

print("Bigrams for transition tags")
train_transition_tag_bigrams = get_ngrams(train_transition_tag_tokens, 2)
with open('../output/tag_bigrams.txt', 'w') as transition_tag_bigrams_output_file:
    transition_tag_bigrams_output_file.write(str(train_transition_tag_bigrams))
print(len(train_transition_tag_bigrams))

```

```

Unigrams for words
30036
Vocab size 30036
Unigrams for word-tag
1126281
41029
Unigrams for tags
470
Unigrams for transition tags
472
Bigrams for transition tags
8255

```

## 4.2

A transition probability is the probability of a tag given its previous tag. Calculate transition probabilities of the training set using the following equation:

$$P(t_i, t_i) = C(t_i, t_i) / C(t_i)$$

```

In [8]: def get_transition_probability(tag_unigrams, tag_bigrams, lambda_value, vocab_size):
        transition_probability = {}
        for i in tag_bigrams:
            previous_tag = i.split(" ")[0]
            transition_probability[i] = (tag_bigrams[i] + lambda_value) / (tag_unigrams[pr
        return transition_probability

```

```
In [9]: print("Get transition probability")
        train_transition_probability = get_transition_probability(train_transition_tag_unigrams)
        # print(train_transition_probability)
```

Get transition probability

### 4.3

An emission probability is the probability of a given word being associated with a given tag. Calculate emission probabilities of the training set using the following equation:

$$P(w_i, t_i) = C(w_i, t_i) / C(t_i)$$

```
In [10]: def get_emission_probability(word_tag_unigrams, tag_unigrams, lambda_value, vocab_size):
        emission_probability = {}
        for i in word_tag_unigrams:
            split_tokens = i.split('/')
            word = split_tokens[0]
            tag = split_tokens[-1]
            emission_probability[i] = (word_tag_unigrams[i] + lambda_value) / (tag_unigrams[tag] + lambda_value)
        # print(word, tag)
        # print(word_tag_unigrams[i], lambda_value)
        # print(tag_unigrams[tag], lambda_value, vocab_size)
        # print(emission_probability[i])
        return emission_probability
```

```
In [11]: print("Get emission probability")
        train_emission_probability = get_emission_probability(train_word_tag_unigrams, train_tag_unigrams, lambda_value, vocab_size)
        # print(train_emission_probability)
```

Get emission probability

4.4 Generate 5 random sentences using the previously learned HMM. Output each sentence (with the POS tags) and its probability of being generated.

Hint: With the help of emission probabilities and transition probabilities collected from 4.2 and 4.3, 1. Start with " tag. 2. Choose next tag based on random choice but considering probabilities e.g. tag draw = random.choices(„). 3. Now choose word for the corresponding tag using emission probabilities (all the words that can be generated from that tag and corresponding probabilities they can be generated with.)e.g. word draw = random.choices(„). 4. Keep repeating steps 2 and 3 till you hit end token '.' 5. Report the sentence and the probability with which this sentence can be generated.

```
In [12]: def generate_sentence(transition_probability, emission_probability):
        sentence_tags = []
        sentence_words = []
        sentence_transition_probability = []
        sentence_emission_probability = []
```

```

while(True):
    # Get potetial next set of tags and its probabilities
    current_tag = "START"
    next_tag = []
    next_tag_probabilities = []
    next_tag_dict = {}
    for i in transition_probability:
        if i.split(" ")[0] == current_tag and i != "END START":
            tag = i.split(" ")[1]
            probability = transition_probability[i]
            next_tag.append(tag)
            next_tag_probabilities.append(probability)
            next_tag_dict[tag] = probability

    # Random pick of a tag
    tag_drawn = random.choices(next_tag, next_tag_probabilities)[0]

#     print("Tag")
#     print(next_tag)
#     print(next_tag_probabilities)
#     print(tag_drawn)

    # Limiting number of words in the sentence to 30
    if (tag_drawn == "END") | (len(sentence_words)==30):
        break

    # Get potential words to be chosen out of a selected tag
    next_word = []
    next_word_probabilities = []
    next_word_dict = {}
    for i in emission_probability:
        if i.split("/")[-1] == tag_drawn:
            word = i.split("/")[-2]
            probability = emission_probability[i]
            next_word.append(word)
            next_word_probabilities.append(probability)
            next_word_dict[word] = probability

    # Random pick of a word
#     print("Word")
#     print(next_word)
#     print(next_word_probabilities)
    word_drawn = random.choices(next_word, next_word_probabilities)[0]

    sentence_tags.append(tag_drawn)
    sentence_words.append(word_drawn)
    sentence_transition_probability.append(next_tag_dict[tag_drawn])
    sentence_emission_probability.append(next_word_dict[word_drawn])

```



```

        current_tag = next_tag

    total_probability = 1
    for i in range(len(sentence_words)):
        total_probability = total_probability * sentence_transition_probability[i] * s

    sentence = sentence_words[0]
    for i in range(1, len(sentence_words)):
        sentence = sentence + " " + sentence_words[i]
    return(sentence_tags, sentence_words, sentence_transition_probability, sentence_em

In [13]: sentence_generation_file = open('../output/generated_sentences.txt', 'w')
    for i in range(1,6):
        print("Generating sentence: {}".format(i))
        sentence_generation_file.write("Sentence: {}\n".format(1))
        sentence_tags, sentence_words, sentence_transition_probability, \
            sentence_emission_probability, total_probability, sentence = generate_sentence(
            train_transition_probability, train_emission_probability
        )
        sentence_generation_file.write("Words: \n {} \n ".format(sentence_words))
        sentence_generation_file.write("Tags: \n {} \n ".format(sentence_tags))
        sentence_generation_file.write("Transition probability \n {} \n ".format(sentence_
        sentence_generation_file.write("Emission probability \n {} \n ".format(sentence_e
        sentence_generation_file.write("Total probability: {} \n \n".format(total_probabi
        sentence_generation_file.write("Sentence: \n {} \n \n".format(sentence))
    sentence_generation_file.close()

Generating sentence: 1
Generating sentence: 2
Generating sentence: 3
Generating sentence: 4
Generating sentence: 5

```

#### 4.5

For each word in the test dataset, derive the most probable POS tag sequence using the Viterbi algorithm; pseudo-code can be found in the textbook <http://web.stanford.edu/~jurafsky/slp3/ed3book.pdf> under Figure 8.5. Viterbi algorithm should be implemented following the pseudocode provided for reference.

Hint: Traversing through back-pointer data structure at the end of algorithm would provide information about the best possible previous tag. So when you are at the second last word of the sentence, calling back-pointer here would give the tag information for the first word in the sentence.

Submit the output in a file exactly with the following format (where each line contains no more than one pair): < sentenceID = 1 > word/tag word/tag .... word/tag < EOS > < sentenceID = 2 > word/tag word/tag .... word/tag < EOS >

```

In [14]: def get_state_transition_matrix(states, transition_probability):
    state_transition_df = pd.DataFrame(0, index=states, columns=states)

```

[illegible]

```
In [15]: def get_word_state_emission_matrix(words, states, emission_probability):
    word_state_emission_df = pd.DataFrame(0, index=words, columns=states)
    for key in emission_probability:
        split_tokens = key.split('/')
        word = split_tokens[0]
        tag = split_tokens[-1]
        word_state_emission_df.loc[word, tag] = emission_probability[key]
    return word_state_emission_df
```

[illegible]



```

for i in range(len(test_sentences)):
    print("Processing sentence {}".format(i))
    viterbi_output_file.write("< sentence ID = {} >\n".format(i+1))
    #     test_sentences[i] = ['Bella', 'wanted', 'to', 'board', 'the', 'bus', 'to', 'Chi
    #     test_sentences[i] = ['John', 'nailed', 'the', 'board', 'over', 'the', 'window']
    result = viterbi_algorithm(test_sentences[i],
                               list(train_tag_unigrams.keys()),
                               state_transition_df,
                               word_state_emission_df
                               )
    for word,tag in result:
        viterbi_output_file.write("{} / {} \n".format(word,tag))
    viterbi_output_file.write("< EOS >\n")
viterbi_output_file.close()

```

```

Processing sentence 0
Converting word: kilowatt-hour to UNK
Converting word: kilowatts to UNK
Converting word: kilowatt to UNK
Converting word: $8 to UNK

```

/Users/karangm/PycharmProjects/pos\_tagging/venv/lib/python3.6/site-packages/ipykernel\_launcher

The current behaviour of 'Series.argmax' is deprecated, use 'idxmax' instead.

The behavior of 'argmax' will be corrected to return the positional maximum in the future. For now, use 'series.values.argmax' or 'np.argmax(np.array(values))' to get the position of the maximum row.

```

Processing sentence 1
Converting word: kilowatt-hour to UNK
Processing sentence 2
Processing sentence 3
Converting word: out-of-pocket to UNK
Processing sentence 4
Converting word: electric-utility to UNK
Processing sentence 5
Processing sentence 6
Processing sentence 7
Converting word: utility-cost to UNK
Processing sentence 8
Processing sentence 9
Converting word: subtype to UNK
Converting word: distributes to UNK
Processing sentence 10
Processing sentence 11

```

Converting word: whereof to UNK  
Converting word: hereunto to UNK  
Converting word: 11th to UNK  
Converting word: sixty-one to UNK  
Converting word: eighty-sixth to UNK  
Processing sentence 12  
Converting word: Resolution to UNK  
Converting word: 22nd to UNK  
Converting word: Maritime to UNK  
Processing sentence 13  
Converting word: intermissions to UNK  
Processing sentence 14  
Processing sentence 15  
Converting word: whereof to UNK  
Converting word: hereunto to UNK  
Converting word: sixty-one to UNK  
Converting word: eighty-sixth to UNK  
Processing sentence 16  
Processing sentence 17  
Converting word: frugality to UNK  
Processing sentence 18  
Processing sentence 19  
Converting word: Pilgrims to UNK  
Processing sentence 20  
Processing sentence 21  
Converting word: Crombie to UNK  
Converting word: Blatz's to UNK  
Processing sentence 22  
Converting word: Crombie to UNK  
Processing sentence 23  
Converting word: Blatz to UNK  
Converting word: Smithtown to UNK  
Processing sentence 24  
Processing sentence 25  
Processing sentence 26  
Converting word: pegboard to UNK  
Processing sentence 27  
Processing sentence 28  
Processing sentence 29  
Processing sentence 30  
Processing sentence 31  
Processing sentence 32  
Converting word: Mattie to UNK  
Converting word: Toonker to UNK  
Converting word: Burkette to UNK  
Converting word: yanking to UNK  
Processing sentence 33  
Processing sentence 34

Processing sentence 35  
Processing sentence 36  
Processing sentence 37  
Converting word: tramp to UNK  
Processing sentence 38  
Converting word: spellbound to UNK  
Processing sentence 39  
Processing sentence 40  
Processing sentence 41  
Processing sentence 42  
Converting word: Juanita to UNK  
Converting word: Lattimer to UNK  
Processing sentence 43  
Converting word: Randolph to UNK  
Converting word: Joel to UNK  
Converting word: replanted to UNK  
Converting word: Annie to UNK  
Processing sentence 44  
Processing sentence 45  
Processing sentence 46  
Processing sentence 47  
Processing sentence 48  
Converting word: Juanita's to UNK  
Processing sentence 49  
Processing sentence 50  
Converting word: 4,585 to UNK  
Converting word: Fisk to UNK  
Processing sentence 51  
Processing sentence 52  
Processing sentence 53  
Converting word: benchmarks to UNK  
Converting word: profundity to UNK  
Processing sentence 54  
Processing sentence 55  
Converting word: libertarian to UNK  
Processing sentence 56  
Converting word: inalienable to UNK  
Processing sentence 57  
Converting word: Avowed to UNK  
Converting word: freethinkers to UNK  
Processing sentence 58  
Converting word: traditionalistic to UNK  
Processing sentence 59  
Converting word: socially-oriented to UNK  
Processing sentence 60  
Processing sentence 61  
Processing sentence 62  
Processing sentence 63

Processing sentence 64  
Converting word: codified to UNK  
Processing sentence 65  
Processing sentence 66  
Converting word: Hesperus to UNK  
Converting word: Lucifer to UNK  
Processing sentence 67  
Processing sentence 68  
Converting word: Hesperus to UNK  
Processing sentence 69  
Processing sentence 70  
Processing sentence 71  
Processing sentence 72  
Processing sentence 73  
Converting word: Warmly to UNK  
Processing sentence 74  
Processing sentence 75  
Converting word: SX-21 to UNK  
Processing sentence 76  
Converting word: plain-clothesmen to UNK  
Processing sentence 77  
Converting word: Thor's to UNK  
Converting word: Antony to UNK  
Converting word: zing to UNK  
Processing sentence 78  
Processing sentence 79  
Processing sentence 80  
Converting word: Ought to UNK  
Converting word: edifying to UNK  
Converting word: Trial to UNK  
Converting word: anti-Semites to UNK  
Converting word: skull-bashings to UNK  
Converting word: gassings to UNK  
Processing sentence 81  
Processing sentence 82  
Converting word: patriots to UNK  
Converting word: terrorizing to UNK  
Converting word: meanest to UNK  
Converting word: pulverizing to UNK  
Processing sentence 83  
Converting word: Trial to UNK  
Converting word: anti-Semitic to UNK  
Converting word: demoralization to UNK  
Processing sentence 84  
Processing sentence 85  
Converting word: Wansee to UNK  
Converting word: Heydrich to UNK  
Processing sentence 86

Converting word: Trial to UNK  
Processing sentence 87  
Converting word: Trial to UNK  
Converting word: anti-Semitism to UNK  
Processing sentence 88  
Converting word: anti-Semitism to UNK  
Converting word: Jew-baiter to UNK  
Processing sentence 89  
Converting word: Heydrich to UNK  
Converting word: Goering to UNK  
Converting word: Solution to UNK  
Converting word: strangulation to UNK  
Converting word: emigration to UNK  
Processing sentence 90  
Converting word: casualties to UNK  
Processing sentence 91  
Converting word: DePugh to UNK  
Converting word: Lauchli to UNK  
Processing sentence 92  
Converting word: Minutemen to UNK  
Processing sentence 93  
Processing sentence 94  
Converting word: Vietnam to UNK  
Processing sentence 95  
Converting word: Albanians to UNK  
Processing sentence 96  
Converting word: Malinovsky to UNK  
Converting word: exalting to UNK  
Processing sentence 97  
Processing sentence 98  
Processing sentence 99  
Converting word: liberating to UNK  
Processing sentence 100  
Converting word: squashed to UNK  
Converting word: suntan to UNK  
Converting word: semi-inflated to UNK  
Processing sentence 101  
Processing sentence 102  
Processing sentence 103  
Processing sentence 104  
Converting word: dirt-catcher to UNK  
Processing sentence 105  
Processing sentence 106  
Converting word: out-of-sight to UNK  
Converting word: out-of-mind to UNK  
Converting word: trek to UNK  
Processing sentence 107  
Converting word: Soignee to UNK



Processing sentence 108  
Processing sentence 109  
Processing sentence 110  
Processing sentence 111  
Converting word: Jannequin's to UNK  
Converting word: tarantara to UNK  
Converting word: rum-tum-tum to UNK  
Converting word: boom-boom-boom to UNK  
Converting word: chansons to UNK  
Converting word: Jannequin to UNK  
Converting word: Lassus to UNK  
Processing sentence 112  
Converting word: Jean-Marie to UNK  
Converting word: LeClair to UNK  
Converting word: Bodin to UNK  
Converting word: Beismortier to UNK  
Converting word: Corrette to UNK  
Converting word: Mondonville to UNK  
Processing sentence 113  
Converting word: forego to UNK  
Processing sentence 114  
Converting word: out-of-the-way to UNK  
Processing sentence 115  
Converting word: dancelike to UNK  
Processing sentence 116  
Converting word: Elegance to UNK  
Processing sentence 117  
Processing sentence 118  
Converting word: Alwise to UNK  
Processing sentence 119  
Processing sentence 120  
Processing sentence 121  
Processing sentence 122  
Converting word: Disapproval to UNK  
Processing sentence 123  
Converting word: full-dress to UNK  
Processing sentence 124  
Processing sentence 125  
Processing sentence 126  
Processing sentence 127  
Processing sentence 128  
Processing sentence 129  
Processing sentence 130  
Converting word: Stacy to UNK  
Converting word: Forbes to UNK  
Processing sentence 131  
Processing sentence 132  
Processing sentence 133

Processing sentence 134  
Converting word: Kimball to UNK  
Converting word: Stacy to UNK  
Processing sentence 135  
Processing sentence 136  
Processing sentence 137  
Converting word: Soak to UNK  
Processing sentence 138  
Converting word: gullet to UNK  
Processing sentence 139  
Converting word: Stacy to UNK  
Converting word: remarry to UNK  
Converting word: Forbes to UNK  
Processing sentence 140  
Processing sentence 141  
Converting word: Methodism to UNK  
Processing sentence 142  
Converting word: Incurably to UNK  
Converting word: devout to UNK  
Converting word: Greenleaf to UNK  
Converting word: Whittier to UNK  
Converting word: 1807-1892 to UNK  
Converting word: plenary to UNK  
Processing sentence 143  
Converting word: Oberlin to UNK  
Processing sentence 144  
Processing sentence 145  
Processing sentence 146  
Converting word: 1811-1884 to UNK  
Converting word: Lyman to UNK  
Converting word: Beecher to UNK  
Processing sentence 147  
Processing sentence 148  
Converting word: anti-slavery to UNK  
Converting word: Finney to UNK  
Converting word: revivals to UNK  
Processing sentence 149