# Report on

CE887-7-AU — Natural Language Engineering

Assignment 1: Probabilities, Regular Expressions & Language Models Aline Villavicencio October 2018

**Student Name**: Atiwat Onsuwan (ao18409)

### Part 1: Tokenization, Part-of-Speech Tagging (30%)

Q1 (20%) Create a program that reads the text from the following website and identify all the types and tokens before and after lowercasing and lemmatization. Use NLTK functions to perform these tasks, from the url reader to the tokenizer and lemmatizer.

https://www.theguardian.com/music/2018/oct/19/while-my-guitar-gently-weeps-beatles-georgeharrison

The output of the program should contain the following information: This text contains types before lemmatization: This text contains tokens before lemmatization: This text contains types after lemmatization: This text contains tokens after lemmatization:

#### Answer.1

```
import re
import nltk
from nltk import word tokenize
from bs4 import BeautifulSoup
import urllib
from urllib import request
from nltk.stem.wordnet import WordNetLemmatizer
from nltk import wordnet
source=urllib.request.urlopen("https://www.theguardian.com/music/2018/oct/19/while-my-
guitar-gently-weeps-beatles-george-harrison").read().decode("utf-8")
soup=BeautifulSoup(source, 'html5lib')
Text = ''
for para in soup.find all('p'):
   Text += (para.text) #Original Text
print(Text) #DISPLAY ORIGINAL TEXT EXTRACT FROM URL
lowertext=Text.lower() # convert MAIN TEXT into lower
# remove punctuation Before LOWERING
mytext tokens nopunct = [word for word in word tokenize(Text)
if re.search("\w", word)]
#Removing Punctuation After LOWRING
mylowertext tokens nopunct = [word for word in word tokenize(lowertext)
if re.search("\w", word)]
print('TYPES BEFORE LOWERING AND LEMMATIZER ')
print(len(set(mytext_tokens_nopunct)))
print('tokens BEFORE LOWER AND LEMMATIZER ')
print(len(mytext tokens nopunct))
```

```
print('TYPES AFTER LOWER ')
print(len(set(mylowertext tokens nopunct)))
print('tokens AFTER ')
print(len(mylowertext tokens nopunct))
#NOW LAMMATIZING
wordnet lemmatizer = WordNetLemmatizer()
nltk tokens = nltk.word tokenize(lowertext)
mylemma tokens nopunct = [word for word in word tokenize(lowertext)
if re.search("\w", word)]
w1 = []
for w in nltk tokens:
   w1.append(wordnet lemmatizer.lemmatize(w))
#print ((w1))
print('TYPES AFTER LEMMATIZER ')
print(len(set(mylemma tokens nopunct)))
print('TOKENS AFTER LEMMATIZER ')
print(len(mylemma tokens nopunct))
```

Acoustic demos of the song, regarded as one of George Harrison's best compositions, to be included in remastered White Album set

**Ben Beaumont-Thomas** 

Fri 19 Oct 2018 12.09 BST

Last modified on Fri 19 Oct 2018 17.30 BST

Three unheard versions of While My Guitar Gently Weeps, regarded by many as George Harrison's greatest contribution to the Beatles, have been released online. The song was written by Harrison in 1968, after he had studied transcendental meditation with the Maharishi Mahesh Yogi in India. There are two completely unheard versions: The Esher Demo is an acoustic version, complete with beautiful multitracked vocals, recorded at Harrison's house in Surrey, in May 1968, in preparation for the studio recordings that began later that month. It is one of 27 demos recorded at the house included on a forthcoming remastered version of The White Album. Another, titled Acoustic Version Take 2, is a raw, shaky but endearing alternative to the solo Harrison demo, played on acoustic guitar and harmonium, that was previously released as part of the Anthology 3 compilation in 1996. At one point Harrison complains, "Yeah, maybe you'll have to give him his own mic" to someone in the studio, as the harmonium drifts into discordance. The final unheard version is a newly remastered mix of the album track. Like the original, it features an electric guitar part from Eric Clapton - he once said he felt the song was about Harrison's dislocation from the group, who didn't share Harrison's zeal for eastern mysticism following their trip to India. Ringo Starr infamously brought a suitcase full of Heinz beans on the trip because he didn't want to eat spicy food, and left after two weeks. The newly remastered White Album is released 9 November, and, as well as the Esher demos, also features 50 mostly unreleased additional recordings made during the studio sessions.

TYPES BEFORE LOWERING AND LEMMATIZER

189

tokens BEFORE LOWER AND LEMMATIZER

314

TYPES AFTER LOWER

181

tokens AFTER

314

TYPES AFTER LEMMATIZER

181

TOKENS AFTER LEMMATIZER

314

Q.2

- a) Assign part-of-speech (POS) tags to all tokens in the text used above. Use one of the implemented POS taggers in NLTK to do this. (5%)
- b )POS taggers do not always assign correct tags to words. Identify tagging errors in the sentences above and briefly explain why these errors may have been caused.

### Answer.2-a

Acoustic demos of the song, regarded as one of George Harrison's best compositions, to be included in remastered White Album set

Ben Beaumont-Thomas

Fri 19 Oct 2018 12.09 BST

Last modified on Fri 19 Oct 2018 17.30 BST

Three unheard versions of While My Guitar Gently Weeps, regarded by many as George Harrison's greatest contribution to the Beatles, have been released online. The song was written by Harrison in 1968, after he had studied transcendental meditation with the Maharishi Mahesh Yogi in India. There are two completely unheard versions: The Esher Demo is an acoustic version, complete with beautiful multitracked vocals, recorded at Harrison's house in Surrey, in May 1968, in preparation for the studio recordings that began later that month. It is one of 27 demos recorded at the house included on a forthcoming remastered version of The White Album. Another, titled Acoustic Version Take 2, is a raw, shaky but endearing alternative to the solo Harrison demo, played on acoustic guitar and harmonium, that was previously released as part of the Anthology 3 compilation in 1996. At one point Harrison complains, "Yeah, maybe you'll have to give him his own mic" to someone in the studio, as the harmonium drifts into discordance. The final unheard version is a newly remastered mix of the album track. Like the original, it features an electric guitar part from Eric Clapton - he once said he felt the song was about Harrison's dislocation from the group, who didn't share Harrison's zeal for eastern

mysticism following their trip to India. Ringo Starr infamously brought a suitcase full of Heinz beans on the trip because he didn't want to eat spicy food, and left after two weeks. The newly remastered White Album is released 9 November, and, as well as the Esher demos, also features 50 mostly unreleased additional recordings made during the studio sessions.

[('Acoustic', 'ADJ'), ('demos', 'NOUN'), ('of', 'ADP'), ('the', 'DET'), ('song', 'NOUN'), (',', '.'), ('regarded', 'VERB'), ('as', 'ADP'), ('one', 'NUM'), ('of', 'ADP'), ('George', 'NOUN'), ('Harrison', 'NOUN'), (''', 'NOUN'), ('s', 'NOUN'), ('best', 'ADJ'), ('compositions', 'NOUN'), (',', '.'), ('to', 'PRT'), ('be', 'VERB'), ('included', 'VERB'), ('in', 'ADP'), ('remastered', 'ADJ'), ('White', 'NOUN'), ('Album', 'NOUN'), ('set', 'VERB'), ('Ben', 'NOUN'), ('Beaumont-Thomas', 'NOUN'), ('Fri', 'NOUN'), ('19', 'NUM'), ('Oct', 'NOUN'), ('2018', 'NUM'), ('12.09', 'NUM'), ('BST', 'NOUN'), ('Last', 'ADJ'), ('modified', 'VERB'), ('on', 'ADP'), ('Fri', 'NOUN'), ('19', 'NUM'), ('Oct', 'NOUN'), ('2018', 'NUM'), ('17.30', 'NUM'), ('BST', 'NOUN'), ('Three', 'NOUN'), ('unheard', 'ADJ'), ('versions', 'NOUN'), ('of', 'ADP'), ('While', 'ADP'), ('My', 'NOUN'), ('Guitar', 'NOUN'), ('Gently', 'NOUN'), ('Weeps', 'NOUN'), (',', '.'), ('regarded', 'VERB'), ('by', 'ADP'), ('many', 'ADJ'), ('as', 'ADP'), ('George', 'NOUN'), ('Harrison', 'NOUN'), ("', 'NOUN'), ('s', 'VERB'), ('greatest', 'ADJ'), ('contribution', 'NOUN'), ('to', 'PRT'), ('the', 'DET'), ('Beatles', 'NOUN'), (',', '.'), ('have', 'VERB'), ('been', 'VERB'), ('released', 'VERB'), ('online.The', 'NOUN'), ('song', 'NOUN'), ('was', 'VERB'), ('written', 'VERB'), ('by', 'ADP'), ('Harrison', 'NOUN'), ('in', 'ADP'), ('1968', 'NUM'), (',', '.'), ('after', 'ADP'), ('he', 'PRON'), ('had', 'VERB'), ('studied', 'VERB'), ('transcendental', 'ADJ'), ('meditation', 'NOUN'), ('with', 'ADP'), ('the', 'DET'), ('Maharishi', 'NOUN'), ('Mahesh', 'NOUN'), ('Yogi', 'NOUN'), ('in', 'ADP'), ('India', 'NOUN'), ('.', '.'), ('There', 'DET'), ('are', 'VERB'), ('two', 'NUM'), ('completely', 'ADV'), ('unheard', 'ADJ'), ('versions', 'NOUN'), (':', '.'), ('The', 'DET'), ('Esher', 'NOUN'), ('Demo', 'NOUN'), ('is', 'VERB'), ('an', 'DET'), ('acoustic', 'ADJ'), ('version', 'NOUN'), (',', '.'), ('complete', 'ADJ'), ('with', 'ADP'), ('beautiful', 'ADJ'), ('multitracked', 'ADJ'), ('vocals', 'NOUN'), (',', '.'), ('recorded', 'VERB'), ('at', 'ADP'), ('Harrison', 'NOUN'), (''', 'NOUN'), ('s', 'NOUN'), ('house', 'NOUN'), ('in', 'ADP'), ('Surrey', 'NOUN'), (',', '.'), ('in', 'ADP'), ('May', 'NOUN'), ('1968', 'NUM'), (',', '.'), ('in', 'ADP'), ('preparation', 'NOUN'), ('for', 'ADP'), ('the', 'DET'), ('studio', 'NOUN'), ('recordings', 'NOUN'), ('that', 'DET'), ('began', 'VERB'), ('later', 'ADV'), ('that', 'ADP'), ('month', 'NOUN'), ('.', '.'), ('It', 'PRON'), ('is', 'VERB'), ('one', 'NUM'), ('of', 'ADP'), ('27', 'NUM'), ('demos', 'NOUN'), ('recorded', 'VERB'), ('at', 'ADP'), ('the', 'DET'), ('house', 'NOUN'), ('included', 'VERB'), ('on', 'ADP'), ('a', 'DET'), ('forthcoming', 'NOUN'), ('remastered', 'ADJ'), ('version', 'NOUN'), ('of', 'ADP'), ('The', 'DET'), ('White', 'NOUN'), ('Album.Another', 'NOUN'), (',', '.'), ('titled', 'VERB'), ('Acoustic', 'ADJ'), ('Version', 'NOUN'), ('Take', 'NOUN'), ('2', 'NUM'), (',', '.'), ('is', 'VERB'), ('a', 'DET'), ('raw', 'ADJ'), (',', '.'), ('shaky', 'ADJ'), ('but', 'CONJ'), ('endearing', 'VERB'), ('alternative', 'ADJ'), ('to', 'PRT'), ('the', 'DET'), ('solo', 'NOUN'), ('Harrison', 'NOUN'), ('demo', 'NOUN'), (',', '.'), ('played', 'VERB'), ('on', 'ADP'), ('acoustic', 'ADJ'), ('guitar', 'NOUN'), ('and', 'CONJ'), ('harmonium', 'NOUN'), (',', '.'), ('that', 'DET'), ('was', 'VERB'), ('previously', 'ADV'), ('released', 'VERB'), ('as', 'ADP'), ('part', 'NOUN'), ('of', 'ADP'), ('the', 'DET'), ('Anthology', 'NOUN'), ('3', 'NUM'), ('compilation', 'NOUN'), ('in', 'ADP'), ('1996', 'NUM'), ('.', '.'), ('At', 'ADP'), ('one', 'NUM'), ('point', 'NOUN'), ('Harrison', 'NOUN'), ('complains', 'VERB'), (',', '.'), ('"', 'NOUN'), ('Yeah', 'NOUN'), (',', '.'), ('maybe', 'ADV'), ('you', 'PRON'), (''', 'VERB'), ('II', 'NOUN'), ('have', 'VERB'), ('to', 'PRT'), ('give', 'VERB'), ('him', 'PRON'), ('his', 'PRON'), ('own', 'ADJ'), ('mic', 'ADJ'), ('"', 'NOUN'), ('to', 'PRT'), ('someone', 'NOUN'), ('in', 'ADP'), ('the', 'DET'), ('studio', 'NOUN'), (',', '.'), ('as', 'ADP'), ('the', 'DET'), ('harmonium', 'NOUN'), ('drifts', 'VERB'), ('into', 'ADP'), ('discordance.The', 'ADJ'), ('final', 'ADJ'), ('unheard', 'ADJ'), ('version', 'NOUN'), ('is', 'VERB'), ('a', 'DET'), ('newly', 'ADV'), ('remastered', 'VERB'), ('mix', 'NOUN'), ('of', 'ADP'), ('the', 'DET'), ('album', 'ADJ'), ('track', 'NOUN'), ('.', '.'), ('Like', 'ADP'), ('the', 'DET'), ('original', 'ADJ'), (',', '.'), ('it', 'PRON'), ('features', 'VERB'), ('an', 'DET'), ('electric', 'ADJ'), ('guitar',

'NOUN'), ('part', 'NOUN'), ('from', 'ADP'), ('Eric', 'NOUN'), ('Clapton', 'NOUN'), ('-', 'NOUN'), ('he', 'PRON'), ('once', 'ADV'), ('said', 'VERB'), ('he', 'PRON'), ('felt', 'VERB'), ('the', 'DET'), ('song', 'NOUN'), ('was', 'VERB'), ('about', 'ADP'), ('Harrison', 'NOUN'), ("', 'NOUN'), ('s', 'ADJ'), ('dislocation', 'NOUN'), ('from', 'ADP'), ('the', 'DET'), ('group', 'NOUN'), (',', '.'), ('who', 'PRON'), ('didn', 'VERB'), ("', 'NOUN'), ('t', 'NOUN'), ('share', 'NOUN'), ('Harrison', 'NOUN'), ("', 'NOUN'), ('s', 'VERB'), ('zeal', 'NOUN'), ('for', 'ADP'), ('eastern', 'ADJ'), ('mysticism', 'NOUN'), ('following', 'VERB'), ('their', 'PRON'), ('trip', 'NOUN'), ('to', 'PRT'), ('India', 'NOUN'), ('.', '.'), ('Ringo', 'NOUN'), ('Starr', 'NOUN'), ('infamously', 'ADV'), ('brought', 'VERB'), ('a', 'DET'), ('suitcase', 'NOUN'), ('full', 'ADJ'), ('of', 'ADP'), ('Heinz', 'NOUN'), ('beans', 'NOUN'), ('on', 'ADP'), ('the', 'DET'), ('trip', 'NOUN'), ('because', 'ADP'), ('he', 'PRON'), ('didn', 'VERB'), (''', 'ADJ'), ('t', 'NOUN'), ('want', 'VERB'), ('to', 'PRT'), ('eat', 'VERB'), ('spicy', 'NOUN'), ('food', 'NOUN'), (',', '.'), ('and', 'CONJ'), ('left', 'VERB'), ('after', 'ADP'), ('two', 'NUM'), ('weeks.The', 'ADJ'), ('newly', 'ADV'), ('remastered', 'VERB'), ('White', 'NOUN'), ('Album', 'NOUN'), ('is', 'VERB'), ('released', 'VERB'), ('9', 'NUM'), ('November', 'NOUN'), (',', '.'), ('and', 'CONJ'), (',', '.'), ('as', 'ADV'), ('well', 'ADV'), ('as', 'ADP'), ('the', 'DET'), ('Esher', 'NOUN'), ('demos', 'NOUN'), (',', '.'), ('also', 'ADV'), ('features', 'VERB'), ('50', 'NUM'), ('mostly', 'ADV'), ('unreleased', 'ADJ'), ('additional', 'ADJ'), ('recordings', 'NOUN'), ('made', 'VERB'), ('during', 'ADP'), ('the', 'DET'), ('studio', 'NOUN'), ('sessions', 'NOUN'), ('.', '.')]

#### Answer.2-b

- 1.('gently', 'NOUN') = Adverb
- 2.('forthcoming','NOUN') = Adjective
- 3.('because','ADP') = conjunction

## Part 2: Regular Expressions, FSAs, and FSTs (30%)

In this part of the assignment you will do some simple information extraction, namely the identification of telephone numbers in text.

Q.3 (20%) Write a regular expression that can find all telephone numbers in a text. Your expression should be able to deal with different formats, for example +55 51 33083838, 1206 872020, 01206 872020 and 05679401945 as well as +44 5679401945 and 0044 5679401945. For full marks: include the output of a Python program that applies your regular expression to any url specified by the user, reads it and finds the telephone numbers. The output should clearly identify what the telephone number is:

Found a match!

Telephone: 01206872020

### Answer.3

```
import nltk
import re
text = (" so I need to match +55 51 33083838, 1206 872020,01206 872020 ,05679401945 ,
+44 5679401945 , 0044 5679401945 Can y" )
```

```
print('Found a matchs')
print('Telephones:')
\label{eq:phone_regex} $$ phone_regex = re.compile("([\+|0|1])\s*(\d{1,10})\s*(\d{1,})\s*(\d{1,})") $$
groups = phone regex.findall(text)
i=1
for g in groups:
   print("".join(g))
```

### Output.3

Found a match!

**Telephones:** 

+555133083838

1206872020

01206872020

05679401945

+445679401945

00445679401945

Q.4 (10%) Write a FSA equivalent to the regular expression you just wrote. You can either use a drawing program or write down a transition table.

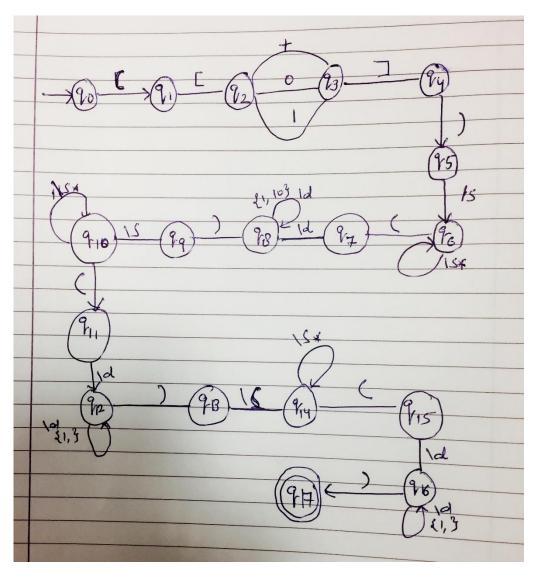
Answer.4

# Regular Expressions: ([\+|0|1])\s\*(\d{1,10})\s\*(\d{1,})\s\*(\d{1,})

# **State Transition Table**

State's	q0	q1	q2	q3	q4	q5	q6	q7	q8	q9	q10	q11	q12	q13	q14	q15	q16	q17
q0		(																
q1			[															
q2				+ 0 1														
q3					]													
q4						)												
q5							\s											
q6							\s*	(										
q7									\d									
q8									{1,10}	)								
q9											\s							
q10											\s*	(						
q11													\d					
q12													{1,}	)				
q13															\s			
q14															\s*	(		
q15																	\d	
q16																	<b>{1, }</b>	)
*q17	L				l	L	L	L		l			L		L			FINAL

### **Finite State Automata**



# Part 3: N-gram models (40%)

Q.5 (10%) Computing a unigram model

- a) Use the Toy dataset. The vocabulary is the words in the sampledata.vocab.txt plus the UNK token. Do not include and in the vocabulary. a) Compute the probabilities in a unigram language model without smoothing. Show your work for P(a), P(c), P(UNK). Create a table in the following format and list all of the probabilities in the unigram model.
- b) Smooth the model using Laplace smoothing. Show your work for P(a), P(c), P(UNK). Show all the smoothed probabilities in a table.

### Program for questions No .5,6 and 7

```
import nltk #import nltk package
import collections # import collections package
from collections import Counter, defaultdict # import counter and defaultdict to count
the number of occurrences
import re # import regular expression package
from tabulate import tabulate #import package 'tabulate' to display table
\# this function reads the files and split words by space (\s+) and define number of
array by \n
def read file(file path):
   with open(file path) as f:
        return [re.split("\s+", line.rstrip('\n')) for line in f]
# this function needs 2 parameters from read file ## to remove <s> and </s> from the
def remove cover(text, cond):
    if cond == "model": #do if cond given = "model"
        for a in (text):
            del a[0]
            del a[-1]
        remove cover.data = []
        for a \overline{in} (text):
            remove cover.data += a
    if cond == "prob":#do if cond given = "prob"
        for a in (text):
            del a[0]
            del a[-1]
# this function is used to calculate UNIGRAM probability of each word and UNK ## 2
parameters are used in this function 1) data set without \langle s \rangle, \langle / s \rangle 2) number of UNK
def unigram model(data set, unk):
   data occurrence = collections.Counter(data set) # 'data occurrence' is stored
number of occurrences of each word in data set
   data_index = list(data_occurrence) # 'data_index' stores list of 'data_occurrence'
    letter count = 0 # define 'letter count' to 0
    unigram prob = [] # set 'unigram prob' as list
    for values in data occurrence: # loop 'data occurrence' to keep values of
'data occurrence' and count number of letters in the data set
       letter count += data occurrence[values]
    i = 0 # set 'i' = 0 and is used to point the index of 'data_index'
    for c in data occurrence: # loop 'data occurrence' to add values into
'unigram_prob'
        unigram_prob.append((data_index[i], round((data_occurrence[c] / letter_count),
5))) # using append to add value to 'unigram prob' and calculate unigram probability
        # using unigram formula == (word occurrence) / number of words in data set
        i += 1
    unigram prob.append(("UNK", round((unk / (letter count + unk)), 5))) # using
append to add value to 'unigram prob' and calculate unigram probability of UNK
    print(tabulate(unigram_prob, headers=['Word', 'Probability'], tablefmt='orgtb1'))
#display unigram probability in table
# this function is used to calculate UNIGRAM Laplace smoothing probability of each
word and UNK
def unigram model laplace(data set, unk, vocab, uncover sentence, is sentence): #and
also calculate UNIGRAM sentence probability
    # 5 parameters are using in this function 1)data set without <s>,</s> 2)number of
UNK found 3) number of vocabulary 4) sentence data set 5) condition for calculate
```

```
unigram sentence probability
   data occurrence = collections.Counter(data set) # 'data occurrence' is stored
number of occurrences of each word in data set
   data index = list(data occurrence) # 'data index' stores list of 'data occurrence'
    letter count = 0 # define 'letter count' to 0
    unigram_laplace_prob = [] # set 'unigram_laplace_prob' as list
    for values in data occurrence: # loop 'data occurrence' to keep values of
'data occurrence' and count number of letters in the data set
       letter_count += data_occurrence[values]
    i = 0 # set 'i' = 0 and is used to point the index of 'data index'
    if is sentence != "true": # do if this function is called with 'is sentence' not
equal to "true"
       for c in data occurrence: # loop 'data occurrence' to add values into
'unigram laplace prob'
            unigram laplace prob.append((data index[i], round((data occurrence[c] + 1)
/ (vocab + letter_count + 1), 5))) # UNIGRAM LAPLACE smoothing formula is used here
            #(word ocurrence +1) / (number of vocabulary + number of words +1)
    if is sentence != "true": # do if this function is called with 'is sentence' not
equal to "true"
        unigram_laplace_prob.append(("UNK", round(((unk + 1) / (vocab + letter_count +
unk + 1)), 5))) # using append to add value to 'unigram_laplace_prob' and calculate
unigram probability of UNK
        print(tabulate(unigram laplace prob, headers=['Word', 'Probability'],
tablefmt='orqtbl')) #display unigram smoothing probability in table
    sentence prob = [[]] * i # define 'sentence prob' as list of list * i to store
values of each word probability
    i = 0
   if is sentence == "true": # do if this function is called with 'is sentence'
equal to "true" (calculate sentence probability)
       unk -= 1
        for c in data occurrence:
           sentence_prob.append([data_index[i], round((data_occurrence[c] + 1) /
(vocab + letter\_count + 1), 5)]) # used the same UNIGRAM LAPLACE calculation above to
store each word probability in 'sentence prob'
           i += 1
        one d array = [item for sublist in sampledata vocab for item in sublist] #
turn sampledata vocab in to 1D array as 'one d array'
       vocab_set = set(one_d_array) # set 'vocab_set' as set of values in
'one d array'
       sentence_set = set(uncover_sentence) # set 'sentence set' as set of values in
'uncover sentence' (sentence data set)
       sentence unk = sentence set.intersection(vocab set) # intersection each of
above set to get UNK in sentence data set
        my list = list(sentence unk) # define 'my list' as list of 'sentence unk'
        j = len(sentence prob) # define 'j' as number of length of sentence prob
        unigram model laplace.num = 0 # 'unigram_model_laplace.num' = 0 (this
parameter will \overline{b}e used to keep value of each sentence probability)
        \mathtt{sum} = \mathtt{0} # set \mathtt{sum} = \mathtt{0} to summation each of the word probability in a sentence
        count = 0 # set count = 0 to count number of loop and will use to check
condition to pass value to 'sum'
        sentence unk = (len(sentence set)) - (len(sentence unk)) # check if there is
UNK in a sentence
        for i in my_list: # outer loop of 'my_list'
            for h in sentence prob: #inner loop of item in 'sentence prob'
                if (sentence prob[j - 1][0]) == i:
                    sum = sentence_prob[j - 1][1] # add every value of word in
sentence to 'sum'
                    count += 1 # now plus 1 to 'count'
                j -= 1 # decrease value of 'j'
                if count == 1:
                    unigram model laplace.num = sum # pass valuse of sum to
```

```
unigram_model_laplace.num
                          if count > 1 and sentence prob[j][0] == i:
                                 unigram model laplace.num = round(unigram model laplace.num * sum,
5) # after 1st loop 'unigram model laplace.num' will start to calculate sentence
probability by multiply every word probability together
                    j = len(sentence_prob) # reset 'j'
             if sentence unk >= 1: # if there is UNK in a sentence
                    unigram_model_laplace.num = round(
                          \label{laplace.num * (((unk + 1) / (vocab + letter\_count + unk + 1)) / (vocab + letter\_count + unk + 1) / (vocab + lett
+ 1)) ** sentence unk), 5) # unk probability will be add and calculate in to the
sentence probability
# this function will check number of UNK in data set compare with vocab
def unk_count(vocab_lenght, data_uncover): # 2 parameters are used in this 1)vocab 2)
dataset with no <s>, </s>
      unk count.list = [] # unk count.list as list
      for a in sampledata vocab: # loop to add every vocab into 'unk count.list'
             if vocab lenght >= 1:
                    unk count.list += a # adding vocab to list
      set_sample_data = set(data_uncover) # set 'set_sample data' as a set of
'data uncover'
      set_vocab_sample_data = set(unk_count.list) # set 'set_vocab_sample_data' as a set
of 'unk count.list'
      unk count.unk = len(
             (set sample data.union(set vocab sample data) -
set sample data.intersection(set vocab sample data))) # now check if we have UNK by
union then - with data intersection vocab
# this function creates bigram pairs of data set
def into bigram pair(data): # 1 parameter is used here is 'data' (original data set
with \langle s \rangle, \langle /s \rangle)
      data = [item for sublist in data for item in sublist] # loop to make item in data
set into one list
      text = " ".join(data) # text = join every item in data using space to split
      into bigram pair.data = list(nltk.bigrams(text.split())) # this called nltk
package to split words into pairs (bigram)
      index = len(into bigram pair.data) - 1 # set 'index' to length of
into bigram pair.data to use for looping
      for item in into bigram pair.data: # loop to remove the pair that contains </s> as
the first word in the pair
             if into bigram pair.data[index][0] == "</s>":
                    del into bigram pair.data[index] # del that pair 'into bigram pair.data'
at 'index' count
             index -= 1 # decrease index evertime of the loop
# this fucntion is use to calculate the probability of BIGRAM unsmoothing using
dataset without cutting starting and ending sentence
def bigram_unsmoothed(data, bigram_data, unmatched_set): #and also uses bigram pairs
of data set, and 'unmatched set' means the posible pairs that not occur in the bigramd
data
      unmatched set = collections.Counter(unmatched set) #count the number of unmatched
pair
      bigram unsmoothed.ind = collections.Counter(data) # count the number of each word
occured in the data set given
      bigram_unsmoothed.pair = collections.Counter(bigram_data) # count the number of
pairs in the data set given
      bigram list = list(bigram unsmoothed.ind) # 1
      pair values = list(bigram unsmoothed.pair) # 2
      unknown pair = list(unmatched set) # 3 # from 1-3 turn the counting information
above into lists
      i = 0
      j = len(bigram list) # count the length of data set counting
      bigram unsmoothed prob = [] # define 'bigram unsmoothed prob' as a list
```

```
for c in bigram unsmoothed.pair: # outer loop = 'bigram unsmoothed.pair' items
        for x in range(j): # inner loop = 'length of data set count'
            if i == len(pair values): # break before number of index in
'pair values' exceeding the last index
                break
            if bigram list[j - 1] == pair values[i][0]: # do if finds word match with
the first pair word
                x = (list(bigram_unsmoothed.ind.values())[j - 1]) # pass values of the
number of pair into 'x'
               bigram unsmoothed_prob.append((pair_values[i],
round((bigram unsmoothed.pair[c]) / (x), 5))) # use the bigram formula to calculate
the prob. of each pair and add the value into 'bigram unsmoothed prob'
                ## along with the particular pair data
                i += 1
            j -= 1
        j = len(bigram list) # reset the length of j to loop again until the outer
loop ends
    j = len(bigram list) # again set j equal to length of counted data set
    i = 0
    for c in unmatched set: # this loop will do the same way of the above loop but use
the not occur pair in the bigram of data set
        for x in range(len(bigram list)):
            if i == len(unknown_pair):
            elif bigram list[j - 1] == unknown pair[i][0]:
                bigram unsmoothed prob.append((unknown pair[i], 0)) # we don't need
formula here because this always 0 (not like add-one smoothing that can have values)
                i += 1
            j -= 1
        j = len(bigram_list)
    i = 0
    j = len(bigram list)
    for c in bigram_unsmoothed.pair: # this loop used to calculate the prob. of pair
that has UNK or UNK and UNK in the pair
        if i == j:
           break # use break before counting of 'i' exceed the last index in the list
in 'bigram list'
       num = (list(bigram unsmoothed.ind.values())[i]) # store value of counted pair
in data set
       bigram unsmoothed prob.append(((bigram list[i], 'UNK'), round(unk count.unk /
(num + len(bigram_list)), 5))) # now add the values of calculation into
'bigram unsmoothed prob' list to use later
       bigram_unsmoothed_prob.append((("UNK", bigram_list[i]), round(unk_count.unk /
len(bigram list), 5))) # now add the values of calculation into
'bigram unsmoothed prob' list to use later
        i += 1
   bigram unsmoothed prob.append((("UNK", "UNK"), unk count.unk / len(bigram list)))
# calculate prob. of UNK and UNK in the same pair and add into
'bigram unsmoothed prob'
    print(tabulate(bigram_unsmoothed_prob, headers=['Pairs occurrence',
'Probability'], tablefmt='orgtbl')) # display all of the prob. in table
# this function does the same method of function 'bigram unsmoothed' but use the
LAPLACE smoothing formula to calculate the prob. of bigram
def bigram smoothed laplace(data, bigram data, unmatched set): # use the same
parameters as 'bigram_unsmoothed' uses
   unknown pair = list(unmatched set)
   unk count.unk += 1
   bigram smoothed laplace.ind = collections.Counter(data)
    list data count = list(bigram smoothed laplace.ind)
   bigram smoothed laplace.pair = collections.Counter(bigram data)
   arr = list(bigram_smoothed_laplace.pair)
```

```
i = 0
    j = len(list data count)
   bigram smoothed laplace.store prob = []
   bigram smoothed laplace prob = []
    for C in bigram smoothed laplace.pair:
        for x in range(len(list data count)):
            if i == len(arr):
                break
            if list data count[j - 1] == arr[i][0]:
                counted data letter = list(bigram smoothed laplace.ind.values())[j -
1]
                counted pair = bigram smoothed laplace.pair[C]
                bigram smoothed laplace.store prob.append([arr[i], ((((counted pair +
1) / (
                        counted data letter + len(list data count)))))))
                bigram smoothed laplace prob.append((arr[i], round(((counted pair + 1)
/ (
                        counted data letter + len(list data count))), 5)))
                i += 1
            j -= 1
        j = len(list_data_count)
    i = 0
    k = 2
    for x in unknown pair:
        if i == len(unknown pair) or k < 0:</pre>
            break
        counted_data_letter = list(bigram_smoothed_laplace.ind.values())[k]
        bigram smoothed laplace.store prob.append([unknown pair[i], ((1 /
((counted data letter) + len(list data count))))])
        bigram smoothed laplace prob.append((unknown pair[i], round(1 /
((counted data letter) + len(list data count)), 5)))
       k = 1
        i += 1
    i = 0
    for c in bigram_smoothed_laplace.pair:
       if i == j:
           break
        counted data letter = list(bigram smoothed laplace.ind.values())[i]
        bigram smoothed laplace.store prob.append(
            [(list_data_count[i], 'UNK'), (unk_count.unk / (counted_data_letter +
len(list data count)))])
        bigram_smoothed_laplace_prob.append(((list_data_count[i], 'UNK'),
round((unk count.unk / (counted data letter + len(list data count))), 5)))
        bigram smoothed laplace.store prob.append([("UNK", list data count[i]),
(unk count.unk / len(list data count))])
        bigram smoothed laplace prob.append((("UNK", list data count[i]),
round(unk_count.unk / len(list_data_count), 5)))
        i += 1
   bigram_smoothed_laplace.store_prob.append([("UNK", "UNK"), (unk_count.unk /
len(list data count))])
    bigram_smoothed_laplace_prob.append((("UNK", "UNK"), round(unk count.unk /
len(list data count), 5)))
    print (tabulate (bigram smoothed laplace prob, headers=['Pairs occurrence',
'Probability'], tablefmt='orgtbl'))
# this function calculates sentence probability by multiplying each of bigram data in
each sentence using values of bigram laplace smoothing we have stored
def bigram sentence probabality(data, text bigram):
   bigram sentence probabality.prob = 0
    i = 0
   count = 0
   num = 0
```

```
j = len(text bigram)
    for c in data:
        for x in text bigram:
            if text bigram[j - 1][0] == data[i]:
                num = text_bigram[j - 1][1] # pass the value of 'text bigram' to 'num'
                count += 1
            j -= 1
            if count == 1:
                bigram sentence probabality.prob = num # pass the value of num to
'bigram sentence probabality.prob'
            if count > 1 and text bigram[j][0] == data[i]:
                bigram sentence probabality.prob = bigram sentence probabality.prob *
num # now multiply them together
        i += 1
        j = len(text_bigram)
# this function calculates all of the possible pair in the BIGRAM
def all possible pair():
    all possible pair.store all pos = []
    flat = [item for sublist in sampledata vocab for item in sublist] # makes the list
of vocab list into just list of vocab flat.insert(0, "<s>") # add <s> at first index
    flat.append("</s>") # add </s> at the last index
    for c in flat: # loop at number of item in 'flat' (new vocab with <s>,</s>) by
'flat'
        for x in flat:
            all possible pair.store all pos.append((c, x)) # add all of the possible
pair into 'store_all_pos'
    all possible pair.store all pos = list(all possible pair.store all pos)
    i = len(all possible pair.store all pos) - 1
    for a in all possible_pair.store_all_pos: # now loop agian to remove </s> that
occurs in the first word of pair
        if all_possible_pair.store_all_pos[i][0] == "</s>":
            del all_possible_pair.store_all_pos[i] # this is how I remove that item
            i = i - 1
    i = len(all_possible_pair.store_all_pos) - 1
   unwanted pair = [('<s>', '<s>'), ('<s>', '</s>')] # set the pairs that we don't
want in the bigram data
    y = 0
    for i in range(26): # loop 26 time (from a-z) to add the pair that we don't to
calculate into 'unwanted pair' means that a-z followed by starting sentence
        unwanted_pair.append((chr(ord('a') + y), '<s>'))
        v += 1
    set unwanted = set(unwanted pair) #1
    prm first set = set(all possible pair.store all pos) #2
   prm second set = set(into bigram pair.data) #4
    unmatched_set = prm_first_set - prm_second_set #5
    all_possible_pair.unmatched_set = unmatched_set - set_unwanted # now I extract and
get only the pair that we don't have yet in the bigram data by substracting each set
    all_possible_pair.sampledata_cover = [item for sublist in read_sampledata_cover
for item in sublist] # turn the 'read_sampledata_cover' (data set with <s>, </s>) into
sampledata cover as list
# this function will check the number of UNK in the sentence test given and change
every word in the sentence that is UNK to UNK
def sort list(sentence data): # 1 parameter needs is sentence data
    sentence = [item for sublist in sentence data for item in sublist] #1
    store raw vocab = [item for sublist in all possible pair.store all pos for item in
sublist] #2
    set sentence = [item for sublist in sentence data for item in sublist] #3
    # from 1-3 I changed every list of list into list to turn them into 'set'
    set_sentence = set(set_sentence) # turn list of sentence into set
```

```
set vocab = set(store raw vocab) # turn list of vocab into set
    intersec bigram and vocab = set sentence - set vocab # now check the UNK in the
sentence by substract the sentence with the set of vocab
   arr sentence = []
    arr intersec bigram and vocab = []
    for s in sentence: # loop to add word in the sentence in to 'arr sentence'
       arr_sentence.append([s])
   for s in intersec bigram and vocab: # loop to add the words that are UNK in
'arr intersec bigram and vocab'
       arr intersec bigram and vocab.append([s])
    k = 0
    for i in arr sentence:
       for j in arr_intersec_bigram_and_vocab:
           if i == j:
               del arr sentence[k] # loop to delete that UNK word in the sentence
               arr sentence.insert(k, ["UNK"]) # and insert "UNK" instead
       k = k + 1
    sort_list.new_sentence = [[]] # define list of list
    for s in arr sentence: # now turn into list of list again with new value
        sort list.new sentence.append([s])
###### read the files by calling 'read file' function and store the values in
parameter as list#####
read sampledata = read file('sampledata.txt') #
read_sampledata_cover = read_file('sampledata.txt')
sampledata vocab = read file('sampledata.vocab.txt')
sentence data = read file('sampletest.txt')
###### read the files by calling 'read file' function and store the values in
parameter as list#####
vocab_lenght = len(sampledata_vocab) # 'vocab_lenght' as the 'sampledata_vocab' length
remove_cover(read_sampledata, "model") # call function 'remove_cover' to remove <s>,
</s> of sample data
unk count (vocab lenght, remove cover.data) # count the number of UNK by using function
'unk count'
print("Training data (sampledata.txt) : ", read sampledata cover) #print the read
files of 'read sampledata cover'
print("Vocabulary (sampledata.vocab.txt) : ", sampledata vocab) #print the read files
of 'sampledata vocab'
print("UNK Found : ", unk count.unk) # display the number of UNK
print("\nxxxxxxxxxxx U N I G R A M xxxxxxxxxxxxx (UNSMOOTHED)")
unigram model (remove cover.data, unk count.unk) #calculate the UNIGRAM prob. of word
using 'unigram model' function
print("xxxxxxxxxxxx U N I G R A M xxxxxxxxxxxxx (SMOOTHED)")
unigram_model_laplace(remove_cover.data, unk_count.unk, len(unk_count.list), 0,
"false") # calculate the UNIGRAM LAPLCE SMOOTHING prob. of word using
'unigram model laplace' function
print("\nxxxxxxxxxxxxxx B I G R A M xxxxxxxxxxxxxx (UNSMOOTHED)")
into bigram pair (read sampledata cover) # send the sample data without removing <s>,
</s> to get pairs of BIGRAM MODEL
all possible pair() # now use the bigram data from above line to calculate all
possible pairs in the bigram prob. table
bigram unsmoothed(all_possible_pair.sampledata_cover, into_bigram_pair.data,
all possible pair.unmatched set)
# send all the information we got through the 'bigram unsmoothed' function to
calculate bigram unsmoothing prob.
print("\nxxxxxxxxxxxxxxx B I G R A M xxxxxxxxxxxxxxx (SMOOTHED)")
```

```
bigram smoothed laplace(all possible pair.sampledata cover, into bigram pair.data,
all possible pair.unmatched set)
# call 'bigram smoothed laplace' function to calculate bigram laplace smoothing using
the same information as unsmoothing used
print("\nxxxxxxxxxxxxxxx S E N T E N C E _ P R O B A B I L I T Y xxxxxxxxxxxxxxx")
print("Sentence Data(sampletest.txt) : ", sentence data) #display the every sentences
data in the sampletest.txt
remove cover(sentence data, "prob") #call 'remove cover' to remove <s>, </s> from
sentence data
sentence = read file('sampletest.txt') #read the files by calling 'read file' function
and store the values in 'sentence' as list
unigram model laplace(remove cover.data, unk count.unk, vocab lenght,
sentence_data[2], "true") #calculate the unigram laplace smoothing of sentence 3 using
'unigram_model laplace'
print("Sentence no.3", sentence[2], "
                                         : Unigram Prob. :",
unigram model laplace.num) # print the result of sentence 3
unigram model laplace (remove cover.data, unk count.unk, vocab lenght,
sentence_data[3], "true") #calculate the unigram laplace smoothing of sentence 4 using
'unigram model laplace'
print("Sentence no.4", sentence[3], ": Unigram Prob. :", unigram_model_laplace.num) #
print the result of sentence 4
unigram_model_laplace(remove_cover.data, unk_count.unk, vocab_lenght,
sentence data[4], "true") #calculate the unigram laplace smoothing of sentence 5 using
'unigram model laplace'
print("Sentence no.5", sentence[4], ": Unigram Prob. :", unigram model laplace.num) #
print the result of sentence 5
sort list(read file('sampletest.txt')) # call the 'sort list' function to get new
sentence information to use in another function
new sentence = [item for sublist in sort list.new sentence for item in sublist] # turn
the information from 'sort list.new sentence' into list named 'new sentence'
into bigram pair(new sentence) #now send the 'new sentence' through 'into bigram pair'
function
bigram_sentence_probabality(into_bigram_pair.data[10:14],
bigram_smoothed_laplace.store_prob) # calculate bigram laplace smoothing prob. of
sentence 3(array index) using function 'bigram_sentence_probabality'
print("Sentence no.3", sentence[2], " : Bigram Prob. :",
round((bigram sentence probabality.prob), 5)) #display the result of sentence 3 prob.
bigram sentence probabality (into bigram pair.data[14:19],
bigram_smoothed_laplace.store_prob) # calculate bigram laplace smoothing prob. of
sentence 4(array index) using function 'bigram sentence probabality'
print("Sentence no.4", sentence[3], ": Bigram Prob. :",
round((bigram sentence probabality.prob), 5)) #display the result of sentence 3 prob.
bigram sentence probabality (into bigram pair.data[19:24],
bigram smoothed laplace.store prob) # calculate bigram laplace smoothing prob. of
sentence 5 (array index) using function 'bigram sentence probabality'
print("Sentence no.5", sentence[4], ": Bigram Prob. :",
round((bigram_sentence_probabality.prob), 5)) #display the result of sentence 3 prob.
```

### Answer.5 a and b

***	XXXXXXX	XXXXX U	NIGR	A M	xxxxxxxxxxxx	(UNSMOOTHED)
1	Word	Prob	bability	1		
1-		+		-		
1	a	I	0.26667	1		
1	b	1	0.33333	1		
1	С	I	0.4	1		
1	UNK	1	0	1		
X	XXXXXXX	XXXXX U	NIGR	AM	xxxxxxxxxxxx	(SMOOTHED)
1	Word	I Droi	bability			
		FIO	Jubilita			
-		+	-			
	a		-	-1		
				I		
	a		0.26316	-   -   -		

Q.6 (20%) Computing a bigram model Use the Toy dataset. The vocabulary is the words in sampledata.vocab.txt, plus the UNK token and symbols. should be included only in the context or history. should not be included in the history but only as the following word. The table below should clarify for you. a) Compute the probabilities in a bigram language model without smoothing. Show your work for P(b|a), P(UNK|), P(UNK|UNK). Create a table in the following format and list all of the probabilities in the model.

b) Smooth the model using Laplace smoothing. Show your work for P(b|a), P(UNK|), P(UNK|UNK). Show all the smoothed probabilities in a table.

### Answer. 6 a and b

xxxxxxxxxxxxxx B   Pairs occurrence				(UNSMOOTHED)
   (' <s>', 'a')</s>	-+	0.66667	-    <b> </b>	
('a', 'a')	i	0.25	Ī	
('a', 'b')	1	0.5	I	
('b', 'b')	1	0.2	I	
('b', 'c')	1	0.6	I	
('c', 'c')	1	0.33333	I	
('c', '')	1	0.33333	T	
('a', 'c')	1	0.5	T	
('c', 'b')	1	0.16667	1	
(' <s>', 'b')</s>	1	0.33333	1	
('c', 'a')	1	0.16667	1	
('b', '')	1	0.2	I	
(' <s>', 'c')</s>	1	0	I	
('a', '')	1	0	I	
('b', 'a')	1	0	I	
(' <s>', 'UNK')</s>	1	0	I	
('UNK', ' <s>')</s>	1	0	I	
('a', 'UNK')	1	0	I	
('UNK', 'a')	1	0	I	
('b', 'UNK')	1	0	I	
('UNK', 'b')	1	0	I	
('c', 'UNK')	1	0	I	
('UNK', 'c')	1	0	I	
('', 'UNK')	1	0	I	
('UNK', '')	1	0	I	
('UNK', 'UNK')	T	0	I	

### XXXXXXXXXXXXXX B I G R A M XXXXXXXXXXXXXXXXX (SMOOTHED)

Pairs occurrence	1	Probability
	+	0.075
(' <s>', 'a')</s>	- !	0.375
('a', 'a')	!	0.22222
('a', 'b')	- 1	0.33333
('b', 'b')	- 1	0.2
('b', 'c')	ı	0.4
('c', 'c')	- 1	0.27273
('c', '')	- 1	0.27273
('a', 'c')	-1	0.33333
('c', 'b')	- 1	0.18182
(' <s>', 'b')</s>	-1	0.25
('c', 'a')	-1	0.18182
('b', '')	-1	0.2
(' <s>', 'c')</s>	1	0.1
('a', '')	- 1	0.11111
('b', 'a')	1	0.125
(' <s>', 'UNK ')</s>	1	0.125
('UNK', ' <s>')</s>	1	0.2
('a', 'UNK ')	1	0.11111
('UNK', 'a')	1	0.2
('b', 'UNK ')	i	0.1
('UNK', 'b')	i	0.2
('c', 'UNK ')	i	0.09091
('UNK', 'c')	i	0.2
('', 'UNK ')		0.125
('UNK', '')		0.125
		0.2
('UNK', 'UNK')	- 1	0.2

Q.7 (10%) Computing sentence probabilities Use the Toy dataset. There are 5 sentences in sampletest.txt. Using the smoothed models above, compute the probability of each sentence. For unigram probability, you should ignore the and symbols.

- a) Show your work for sentence numbers 3, 4, 5, for each model: unigram and bigram.
- b) Fill in the probabilities of all the sentences in a table.

### Answer, 7 a and b

```
********** S E N T E N C E _ P R O B A B I L I T Y ***************
Sentence no.3 ['<s>', 'c', 'b', 'a', '</s>'] : Unigram Prob. : 0.03062
Sentence no.4 ['<s>', 'a', 'b', 'c', 'd', '</s>'] : Unigram Prob. : 0.00161
Sentence no.5 ['<s>', 'a', 'd', 'e', 'b', '</s>'] : Unigram Prob. : 0.00023
Sentence no.3 ['<s>', 'c', 'b', 'a', '</s>'] : Bigram Prob. : 0.00025
Sentence no.4 ['<s>', 'a', 'b', 'c', 'd', '</s>'] : Bigram Prob. : 0.00091
Sentence no.5 ['<s>', 'a', 'd', 'e', 'b', '</s>'] : Bigram Prob. : 0.00033
```