

# CS 505 Homework 03: N-Gram Modelling

Due Monday 10/9 at midnight (1 minute after 11:59 pm) in Gradescope (with a grace period of 6 hours)

You may submit the homework up to 24 hours late (with the same grace period) for a penalty of 10%.

All homeworks will be scored with a maximum of 100 points; point values are given for individual problems, and if parts of problems do not have point values given, they will be counted equally toward the total for that problem.

Note: I strongly recommend you work in **Google Colab** (the free version) to complete homeworks in this class; in addition to (probably) being faster than your laptop, all the necessary libraries will already be available to you, and you don't have to hassle with `conda`, `pip`, etc. and resolving problems when the install doesn't work. But it is up to you! You should go through the necessary tutorials listed on the web site concerning Colab and storing files on a Google Drive. And of course, Dr. Google is always ready to help you resolve your problems.

I will post a "walk-through" video ASAP on my [Youtube Channel](#).

## Submission Instructions

You must complete the homework by editing **this notebook** and submitting the following two files in Gradescope by the due date and time:

- A file `HW03.ipynb` (be sure to select `Kernel -> Restart and Run All` before you submit, to make sure everything works); and
- A file `HW03.pdf` created from the previous.

For best results obtaining a clean PDF file on the Mac, select `File -> Print Review` from the Jupyter window, then choose `File-> Print` in your browser and then `Save as PDF`. Something similar should be possible on a Windows machine -- just make sure it is readable and no cell contents have been cut off. Make it easy to grade!

The date and time of your submission is the last file you submitted, so if your IPYNB file is submitted on time, but your PDF is late, then your submission is late.

## Collaborators (5 pts)

Describe briefly but precisely

1. Any persons you discussed this homework with and the nature of the discussion;
2. Any online resources you consulted and what information you got from those resources; and
3. Any AI agents (such as chatGPT or CoPilot) or other applications you used to complete the homework, and the nature of the help you received.

A few brief sentences is all that I am looking for here.

www.medium.com : to understand about perplexity

chatGPT : got stuck at perplexity and my solution was giving me wrong outputs. So asked GPT for help in debugging this part

www.towardsdatascience.com : knowing about N gram model and making my head clear

```
In [51]: import math
import numpy as np
from numpy.random import shuffle, seed, choice
import nltk
from tqdm import tqdm
from collections import defaultdict

# First time you will need to download the corpus:
# Run the following and download the book collection

from nltk.corpus import brown
nltk.download('brown')
```

```
[nltk_data] Downloading package brown to
[nltk_data] /Users/mohanthota/nltk_data...
[nltk_data] Package brown is already up-to-date!
```

```
Out[51]: True
```

## Problem One: Bag of N-Grams

A BOW is a language modelling technique (also called a Term Frequency Vector) which creates a frequency distribution for a set of tokens -- or unigrams! Extending this idea a bit, we can also create a Bag of N-Grams, which is a frequency distribution for a set of N-grams for some N. If we divide the frequency by the number of N-grams, we have a probability distribution, such as we showed for the exciting text about John and Mary in Lecture 5.

For this homework, we are going to create such Bag of N-Gram models for  $N = 1, 2, 3,$  & 4, for the sentences in `brown.sents()`. We will evaluate them using a test set, and then in the second part of the homework, we shall use them to generate sentences.

**Note 1:** We do not want to do the same low-level transformations in this project as we did in HW 02. We will keep the capitalization, punctuation, and words in all their various forms. There are some strange things in `brown.sents()`, such as double semicolons, and bad sentence segmentation, but we will assume the processing of the texts into `brown.sents()` was consistent, and we will see what our model makes of this data.

**Note 2:** Since `brown.sents()` contains punctuation marks as well as words, we shall use the term **tokens** for the strings stored in the sentence lists.

## Part A: Randomize the list of sentences and split into training and testing sets

We will use `brown.sents()` (a list of list of tokens) as the basis of our N-gram models. The list `brown.words()` is simply the concatenation of all these lists of tokens.

We will shuffle the list into a random order, but using a seed value so that the order of the random shuffle is the same each time.

1. Read about `numpy.random.seed` and `numpy.random.shuffle`.
2. Set the seed to `0` and shuffle the list: you can't shuffle `brown.sents()` and because `numpy.random.shuffle` modifies the list **in place**, to avoid reshuffling:
  - Convert **`brown.sents()`** to a list and assign to a new variable **`sentences`** and then
  - **Copy `sentences`** to a new variable **`shuffled_sentences`** and then
  - Shuffle that list.

In this way, you will have the original list, and a randomized list, but because of `seed(0)` it will be in the same order every time you run your code (and when we grade it).

1. Then split `shuffled_sentences` into sets `train_sentences` (first 99.9% of the sentences) and `test_sentences` (last 0.1%).
2. Print out the length of the training and testing sets.
3. Print out the first sentence in each of these sets.

In 4 and 5, label the outputs so we know which is which. (Always make outputs easy to understand!)

NOTE: The terms "training set" and "testing set" are very standard, even though we store these in lists (it is possible that there are duplicate sentences).

```
In [52]: # First, shuffle the set of sentences

# your code here

import numpy as np
from nltk.corpus import brown
sentences = list(brown.sents())
np.random.seed(0)
np.random.shuffle(sentences)
train_size = int(0.999 * len(sentences)) # as said used 99% for training
train_sentences, test_sentences = sentences[:train_size], sentences[train_size:]
print("Length of Training Set:", len(train_sentences))
print("Length of Testing Set:", len(test_sentences))
print("\nStart of training set:\n", train_sentences[0])
print("\nStart of testing set:\n", test_sentences[0])
```

Length of Training Set: 57282

Length of Testing Set: 58

Start of training set:

['Muscle', 'weakness', 'did', 'not', 'improve', ',', 'and', 'the', 'patient', 'needed', 'first', 'a', 'cane', ',', 'then', 'crutches', '.']

Start of testing set:

['It', 'is', 'at', 'least', 'as', 'important', 'as', 'the', 'more', 'dramatic', 'attempts', 'to', 'break', 'down', 'barriers', 'of', 'inequality', 'in', 'the', 'South', '.']

## Part B

Now, you must add the beginning `<s>` and ending `</s>` markers to each sentence in both the training and testing lists. Do not make any other changes to the sentences -- you will see that punctuation has been left in, such as periods at the end of sentences. Again, we will see what our models make of this data set.

Print out the first sentence in each of the training and testing sets to check that all is well.

```
In [53]: # put `<s>` at beginning and `</s>` at end of all sentences.
def bracket_sentence(sent):
    return ['<s>'] + sent + ['</s>']

# your code here
train_sentences, test_sentences = ([bracket_sentence(sent) for sent in train_sentences],
                                   [bracket_sentence(sent) for sent in test_sentences])
print("Start of training set:\n", train_sentences[0], "\n\n")
print("Start of testing set:\n", test_sentences[0])
```

Start of training set:

```
['<s>', 'Muscle', 'weakness', 'did', 'not', 'improve', ',', 'and', 'the', 'patient', 'needed', 'first', 'a', 'cane', ',', 'then', 'crutches', '.', '</s>']
```

Start of testing set:

```
['<s>', 'It', 'is', 'at', 'least', 'as', 'important', 'as', 'the', 'more', 'dramatic', 'attempts', 'to', 'break', 'down', 'barriers', 'of', 'inequality', 'in', 'the', 'South', '.', '</s>']
```

## Part C

Complete the following template for a function to extract N-grams from one sentence, and test it for N = 1,2,3,4 for the first sentences in the training set.

```
In [54]: # Return a list of the N-grams for all sentences s

# Store all N-grams as tuples, so that a unigram is (w,), a bigram is (w1,w2)

def get_Ngrams_for_sentence(N,s):
    if N == 0 or N > len(s):
        return []
    return list(zip(*[s[i:] for i in range(N)]))

# your code here
for N in range(1, 5):
    print(get_Ngrams_for_sentence(N, train_sentences[0]), "\n")
```

```
[('<s>',), ('Muscle',), ('weakness',), ('did',), ('not',), ('improve',),
(' ',), ('and',), ('the',), ('patient',), ('needed',), ('first',), ('a',),
('cane',), (' ',), ('then',), ('crutches',), ('.',), ('</s>',)]
```

```
[('<s>', 'Muscle'), ('Muscle', 'weakness'), ('weakness', 'did'), ('did', 'not'),
('not', 'improve'), ('improve', ' '), (' ', 'and'), ('and', 'the'), ('the', 'patient'),
('patient', 'needed'), ('needed', 'first'), ('first', 'a'), ('a', 'cane'),
('cane', ' '), (' ', 'then'), ('then', 'crutches'), ('crutches', ' '),
(' ', ' '), ('.', ' '), ('</s>',)]
```

```
[('<s>', 'Muscle', 'weakness'), ('Muscle', 'weakness', 'did'), ('weakness',
'did', 'not'), ('did', 'not', 'improve'), ('not', 'improve', ' '), ('improve',
' ', 'and'), (' ', 'and', 'the'), ('and', 'the', 'patient'), ('the', 'patient',
'needed'), ('patient', 'needed', 'first'), ('needed', 'first', 'a'),
('first', 'a', 'cane'), ('a', 'cane', ' '), ('cane', ' ', 'then'), (' ', 'then',
'crutches'), ('then', 'crutches', '.'), ('crutches', '.', '</s>')]
```

```
[('<s>', 'Muscle', 'weakness', 'did'), ('Muscle', 'weakness', 'did', 'not'),
('weakness', 'did', 'not', 'improve'), ('did', 'not', 'improve', ' '), ('not',
'improve', ' ', 'and'), ('improve', ' ', 'and', 'the'), (' ', 'and', 'the',
'patient'), ('and', 'the', 'patient', 'needed'), ('the', 'patient', 'needed',
'first'), ('patient', 'needed', 'first', 'a'), ('needed', 'first', 'a', 'cane'),
('first', 'a', 'cane', ' '), ('a', 'cane', ' ', 'then'), ('cane', ' ', 'then',
'crutches'), (' ', 'then', 'crutches', '.'), ('then', 'crutches', '.', '</s>')]
```

## Part D

Now create lists of N-grams for all the sentences in your training set (NOT the testing set). Complete the following template to assign these to the given list.

Print out the number of N-grams, and the first 5 N-grams in each list for N = 1, 2, 3, 4.

Note that this number is the number of occurrences of N-grams, which may not be unique in the list.

```
In [55]: Ngrams = [None]*5      # first slot is empty, then Ngram[1] will hold unigrams
        # your code here

        for N in range(1, 5):
            Ngrams[N] = [ngram for sent in train_sentences for ngram in get_Ngrams_f
            print(f"For N={N}:")
            print("Number of N-grams:", len(Ngrams[N]))
            print("First 5 N-grams:", Ngrams[N][:5])
            print("\n")
```

```

For N=1:
Number of N-grams: 1274667
First 5 N-grams: [('<s>',), ('Muscle',), ('weakness',), ('did',), ('not',)]

For N=2:
Number of N-grams: 1217385
First 5 N-grams: [('<s>', 'Muscle'), ('Muscle', 'weakness'), ('weakness', 'did'), ('did', 'not'), ('not', 'improve')]

For N=3:
Number of N-grams: 1160103
First 5 N-grams: [('<s>', 'Muscle', 'weakness'), ('Muscle', 'weakness', 'did'), ('weakness', 'did', 'not'), ('did', 'not', 'improve'), ('not', 'improve', ',')]

For N=4:
Number of N-grams: 1102821
First 5 N-grams: [('<s>', 'Muscle', 'weakness', 'did'), ('Muscle', 'weakness', 'did', 'not'), ('weakness', 'did', 'not', 'improve'), ('did', 'not', 'improve', ','), ('not', 'improve', ',,', 'and')]

```

## Part E

We will now create a probability distribution for each of the Ngram collections. Note carefully that you must divide the frequency of each N-gram by the number of occurrences of N-grams, not the number of unique N-grams.

Note that we have set the probability of the unigram `<s>` to 1.0. This is because we are interested in calculating the probability of sentences, which *always* begin with the token `'<s>'`.

Complete the following template and then

1. Print out the total number of N-grams in each dictionary (they should be a bit smaller than the totals in the last part - why?).
2. Test your code by printing out the probability of the following Ngrams to 8 digits of precision:

```

('to',)
('to','the')
('to','the','house')
('to','the','house','.')

```



In [56]: *# Create a defaultdict with the frequency distribution for the training set*

```
from collections import Counter
def get_Ngram_distribution(N, Ngrams):
    freq_dist = defaultdict(float)
    for ngram in Ngrams:
        freq_dist[ngram] += 1
    total_ngrams = len(Ngrams)
    for ngram in freq_dist:
        freq_dist[ngram] /= total_ngrams

    return freq_dist

# your code here

# now create for N = 1,2,3,4

Ngram_distribution = [None]*5

# your code here

for N in range(1, 5):
    Ngram_distribution[N] = get_Ngram_distribution(N, Ngrams[N])
Ngram_distribution[1][('<s>',)] = 1.0

# tests

print(f"\nThe Probability of ('<s>',) is {np.around(Ngram_distribution[1][('<s>',)]
print(f"\nThe Probability of ('to',) is {np.around(Ngram_distribution[1][('t
print(f"\nThe Probability of ('to','the') is {np.around(Ngram_distribution[2
print(f"\nThe Probability of ('to','the','house') is {np.around(Ngram_distri
print(f"\nThe Probability of ('to','the','house','.') is {np.around(Ngram_d
```

The Probability of ('<s>',) is 1.0.

The Probability of ('to',) is 0.02017703.

The Probability of ('to','the') is 0.00281341.

The Probability of ('to','the','house') is 9.48e-06.

The Probability of ('to','the','house','.') is 2.72e-06.

## Probability and Perplexity

Now we will calculate the probability and the perplexity of sequences of tokens, using the principle of "Stupid Backoff" as explained in the paper:

<https://aclanthology.org/D07-1090.pdf>

and explicated in this StackOverflow post:

<https://stackoverflow.com/questions/16383194/stupid-backoff-implementation-clarification>

Before describing "Stupid Backoff," let us consider the naive way to calculate the probability of a sequence of tokens which starts with `<s>`.

## A simple and naive way to calculate probabilities of sequences of tokens

Suppose we have a quadrigram model ( $N = 4$ ), we have a sequence of tokens

$$w_1, w_2, \dots, w_n,$$

Assume we have calculated all the N-gram probabilities in `Ngram_distribution[N]` for  $N = 1, 2, 3, 4$ .

We will calculate the probability of each successive token  $w_i$  in as much left context as we have, up to 3 (the last token in the 4-gram being  $w_i$ ).

$$\begin{aligned} p_1 &= \text{Ngram\_distribution}[1][(w_1,)] \\ p_2 &= \text{Ngram\_distribution}[2][(w_1, w_2)] / \text{Ngram\_distribution}[1][(w_1,)] \\ p_3 &= \text{Ngram\_distribution}[3][(w_1, w_2, w_3)] / \text{Ngram\_distribution}[2][(w_1, w_2)] \\ p_4 &= \text{Ngram\_distribution}[4][(w_1, w_2, w_3, w_4)] / \text{Ngram\_distribution}[3][(w_1, w_2, w_3)] \\ &\dots \\ p_i &= \text{Ngram\_distribution}[4][(w_{i-3}, w_{i-2}, w_{i-1}, w_i)] / \text{Ngram\_distribution}[3][(w_{i-3}, w_{i-2}, w_{i-1})] \\ &\dots \\ p_n &= \text{Ngram\_distribution}[4][(w_{n-3}, w_{n-2}, w_{n-1}, w_n)] / \text{Ngram\_distribution}[3][(w_{n-3}, w_{n-2}, w_{n-1})] \end{aligned}$$

Note that if  $w_1$  is `<s>`, then  $p_1 = 1.0$ .

Finally, let

$$P(w_1, w_2, \dots, w_n) = p_1 * p_2 * \dots * p_n.$$

One messy detail is dealing with the possibility of 0 counts (if that is possible); thus, if the numerator in the above expressions is 0, then the entire product is 0 (you want to avoid a divide by 0 error in the denominator).

## What could possibly go wrong?

Well, if our sentence is from our training set, nothing! All the probabilities will have been calculated for all the possible N-grams.

However, when we have a separate training set, we have to account for the fact that **some N-grams (and even some tokens) may occur in the testing set which do not occur in the training set, and so their probability will be 0.** However, we want to make the best estimate of the probability we can!

There are various solutions, which we discussed in lectures 5 and 6, but the simplest (and very effective for large data sets) is "Stupid Backoff," recursively defined as follows for bigrams, trigrams, and quadrigrams, and **using the probability calculations shown above.**

```

PN_stupid_backoff(w1) = p1    as defined above          # if
this is not 0, else:
                        = (frequency of w1 in whole corpus /
number of tokens in whole corpus)

PN_stupid_backoff(w1, w2) = p2    as defined above      #
if this is not 0, else:
                        = 0.4 * P_stupid_backoff(w2)    #
recursive, use previous definition

PN_stupid_backoff(w1, w2, w3) = p3    as defined above
# if this is not 0, else:
                        = 0.4 * P_stupid_backoff(w2,
w3) # recursive, use definition above

PN_stupid_backoff(w1, w2, w3, w4) = p4    as defined above
# if this is not 0, else:
                        = 0.4 *
P_stupid_backoff(w2, w3, w4)    # recursive, previous
definition

```

This accounts for how to backoff when trying to find the probability of some  $w_i$  in a left context of 1, 2, or 3 tokens.

Then we use these calculations instead of  $p_1, p_2$  as in the previous algorithm, for trigrams it would be:

```

P_stupid_backoff(w1, w2, w3, w4, ..., wn) =
PN_stupid_backoff(w1)
*
PN_stupid_backoff(w1, w2)
*

```

```

PN_stupid_backoff(w1, w2, w3)
                                     *
PN_stupid_backoff(w2, w3, w4)
                                     *
                                     ...
PN_stupid_backoff(w(n-2), w(n-1), wn)
                                     *
```

The "discount factor" 0.4 was proposed by the originators of the method, and seems to work well in practice.

This calculation is unnecessary in generative models, since then we will train on the entire corpus, and only use available N-grams to produce sentences.

## Problem 2

Now we will calculate the probability of a sequence of tokens. We will warm up by considering the simple case, and then consider the more complex case, where "stupid backoff" will be used.

### Part A

For this part, complete the following template to create a function which will calculate the probability of a sequence of tokens.

**Since you may want to use this in the stupid backoff version, you should make sure that if the numerator in the conditional probability calculation is 0, immediately return 0, so that there is no possibility of divide by 0 in the denominator.**

Tests are provided following the cell in which you will write your code.

```
In [57]: # Probability of a list of tokens using N-grams
# W a list of tokens

# Calculate the probability of the last token in W, as we did in the calcula
# check numerator: if it is 0, return 0 immediately and do not do the divisi
# division by 0)

# N == len(W)
def PN(N, W):
    ngram = tuple(W)
    numerator = Ngram_distribution[N].get(ngram, 0)
    if len(W) != N or numerator == 0:
        return 0
    if N == 1:
        return numerator
    denominator = Ngram_distribution[N-1].get(ngram[:-1], 0)
    return numerator / denominator if denominator else 0

def P(N, W):
    if len(W) <= N:
        return PN(len(W), W)
    return np.prod([PN(min(i+1, N), W[max(i-N+1, 0):i+1]) for i in range(len(W)-N+1)])
```

The following are tests to make sure your code is working properly. The values printed should be the same.

```
In [58]: # Sentence: He frowned.
# Using a bigram model

a = Ngram_distribution[2][('<s>', 'He')]
b = Ngram_distribution[2][('He', 'frowned')] / Ngram_distribution[1][('He',)]
c = Ngram_distribution[2][('frowned', '.')] / Ngram_distribution[1][('frowned')]
d = Ngram_distribution[2][('.', '</s>')] / Ngram_distribution[1][('.',)]

print('a=', a)
print(' ', PN(2,('<s>', 'He')))
print('b =', b)
print(' ', PN(2,('He', 'frowned')))
print('c=', c)
print(' ', PN(2,('frowned', '.')))
print('d=', d)
print(' ', PN(2,('.', '</s>')))

print('\na*b*c*d: ', a*b*c*d)
print('P(2,...): ', P(2,('<s>', 'He', 'frowned', '.', '</s>')))
```

```

a= 0.002346012148991486
   0.002346012148991486
b = 0.000351478118528877
   0.000351478118528877
c= 0.39264499316157175
   0.39264499316157175
d= 1.0470533150975245
   1.0470533150975245

a*b*c*d:    3.389982137368031e-07
P(2,...):   3.389982137368031e-07

```

```

In [59]: # Sentence: He frowned.
          # Using a trigram model

a = Ngram_distribution[2][('<s>', 'He')]
b = Ngram_distribution[3][('<s>', 'He', 'frowned')] / Ngram_distribution[2][('<s>', 'He')]
c = Ngram_distribution[3][('He', 'frowned', '.')] / Ngram_distribution[2][('He', 'frowned')]
d = Ngram_distribution[3][('frowned', '.', '</s>')] / Ngram_distribution[2][('frowned', '.')]

print('a*b*c*d: ', a*b*c*d)
print('P(3,...): ', P(3, ('<s>', 'He', 'frowned', '.', '</s>')))

a*b*c*d:    9.492186072583041e-07
P(3,...):   9.492186072583041e-07

```

```

In [60]: # Sentence: He frowned.
          # Using a quadrigram model

a = Ngram_distribution[2][('<s>', 'He')]
b = Ngram_distribution[3][('<s>', 'He', 'frowned')] / Ngram_distribution[2][('<s>', 'He')]
c = Ngram_distribution[4][('<s>', 'He', 'frowned', '.')] / Ngram_distribution[3][('<s>', 'He', 'frowned')]
d = Ngram_distribution[4][('He', 'frowned', '.', '</s>')] / Ngram_distribution[3][('He', 'frowned', '.')]

print('a*b*c*d: ', a*b*c*d)
print('P(4,...): ', P(4, ('<s>', 'He', 'frowned', '.', '</s>')))

a*b*c*d:    9.538640808692934e-07
P(4,...):   9.538640808692934e-07

```

```
In [61]: # Sentence: I love NLP
# using trigrams

# This shows that you should test the numerator to see if it is 0, because t
# you can return 0 immediately, and not have the possibility of division by
# which would occur in the cases d below (c is 0 but would not cause divisio

a = Ngram_distribution[2][('<s>', 'I')]
b = Ngram_distribution[3][('<s>', 'I', 'love')] / Ngram_distribution[2][('<s>
c = Ngram_distribution[3][('I', 'love', 'NLP')] / Ngram_distribution[2][('I',
d1 = Ngram_distribution[3][('love', 'NLP', '</s>')]
d2 = Ngram_distribution[2][('love', 'NLP')]

print('a=', a, '\nb=', b, '\nc=', c, '\nd1=', d1, '\nd2=', d2, '\n')

print('a*b*d*d1: ', a*b*c*d1)
print('P(3,...): ', P(3, ('<s>', 'I', 'love', 'NLP', '</s>')))
```

```
a= 0.001128648701930778
b= 0.0015274769289326804
c= 0.0
d1= 0.0
d2= 0.0

a*b*d*d1: 0.0
P(3,...): 0.0
```

## Part B

Now we will develop the probability for a sequence with the possibility that some N-grams, or even some tokens, are not in the training set. We will use the idea of "stupid backoff" explained in lecture.

Complete the following template and verify that it passes all the tests.

```

In [62]: # Probability with stupid backoff
# same as previous, but have to use recursive (or iterative) method instead
# calling Ngram_distribution directly

# W a list of tokens

num_all_tokens = len(brown.words())

# This returns backed-off probability for single N-gram
# len(W) must be N, this will try whole N-gram, then last N-1 tokens, then N

# Assumes W is a tuple

# calculate for a particular length N-gram
# must have N == len(W)

def PN_with_stupid_backoff(N, W):
    ngram = tuple(W)
    if N == 1:
        return Ngram_distribution[1].get(ngram, W.count(ngram[0]) / num_all_
    if ngram in Ngram_distribution[N]:
        denominator = Ngram_distribution[N-1].get(ngram[:-1], 0)
        return Ngram_distribution[N][ngram] / denominator if denominator else
    return 0.4 * PN_with_stupid_backoff(N-1, W[1:])

import numpy as np
def P_stupid_backoff(N, W):
    if len(W) <= N:
        return PN_with_stupid_backoff(len(W), tuple(W))
    return np.prod([PN_with_stupid_backoff(min(i+1, N), tuple(W[max(i-N+1, 0

```

The following are tests to make sure your code is working properly. The values printed should be the same.

```

In [63]: # 'grandstand' is in testing set but not in the training set
# This uses bigrams

P(2, ('<s>', 'where', 'is', 'the', 'grandstand'))

```

Out[63]: 0.0



```
In [64]: a = PN(2, ('<s>', 'where'))
b = PN(2, ('where', 'is'))
c = PN(2, ('is', 'the'))
d2 = PN(2, ('the', 'grandstand'))      # this is 0, so use less context
d1 = PN(1, ('grandstand'))              # this is 0, so use 0.4*d instead of d2
d = list(brown.words()).count('grandstand') / len(brown.words())

print('a =', a, '\nb =', b, '\nc =', c, '\nd2=', d2, '\nd1=', d1, '\nd =', d, '\n')
print(a*b*c*(0.4*d))
print(P_stupid_backoff(2, ('<s>', 'where', 'is', 'the', 'grandstand')))
```

```
a = 1.6428656505542618e-06
b = 0.006166391726133832
c = 0.08182229363508187
d2= 0
d1= 0
d = 8.611840246918683e-07

2.855359302895301e-16
2.855359302895301e-16
```

```
In [65]: # 'grandstand' is in testing set but not in the training set
# This uses trigrams

P(3, ('<s>', 'where', 'is', 'the', 'grandstand'))
```

```
Out[65]: 0.0
```

```
In [66]: a = PN(2, ('<s>', 'where'))      # <= this works
b3 = PN(3, ('<s>', 'where', 'is'))      # this is 0, so try less context
b2 = PN(2, ('where', 'is'))              # <= this works
c = PN(3, ('where', 'is', 'the'))        # <= this works

d3 = PN(3, ('is', 'the', 'grandstand'))  # this is 0, so try less context
d2 = PN(2, ('the', 'grandstand'))         # this is 0, so try less context
d1 = PN(1, ('grandstand'))                # this is 0, so use probability in w
d = list(brown.words()).count('grandstand') / len(brown.words())      # <= t

print('a =', a, '\nb3=', b3, '\nb2=', b2, '\nc =', c, '\nd3=', d3, '\nd2=', d2, '\nd1=',

# every time we try less context, must multiply by 0.4, so we
# use 0.4*b2 instead of b3 and 0.4*0.4*d instead of d3
print(a*(0.4*b2)*c*(0.4*0.4*d))
print(P_stupid_backoff(3, ('<s>', 'where', 'is', 'the', 'grandstand')))
```

```

a = 1.6428656505542618e-06
b3= 0
b2= 0.006166391726133832
c = 0.4197506600707006
d3= 0
d2= 0
d1= 0
d = 8.611840246918683e-07

```

```

2.3436917228933544e-16
2.3436917228933544e-16

```

```

In [67]: a = PN(2, ('<s>', 'The')) # <= this works

b3 = PN(3, ('<s>', 'The', 'grandstand')) # this is 0, so try less context
b2 = PN(2, ('The', 'grandstand')) # this is 0, so try less context
b1 = PN(1, ('grandstand',)) # this is 0, so have to use frequ
b = list(brown.words()).count('grandstand') / len(brown.words()) # <=

c3 = PN(3, ('The', 'grandstand', 'fell')) # this is 0, so try less context
c2 = PN(2, ('grandstand', 'fell')) # this is 0, so try less context
c1 = PN(1, ('fell',)) # ok, this works

d3 = PN(3, ('grandstand', 'fell', 'down')) # this is 0, so try less context
d2 = PN(2, ('fell', 'down')) # <= this works

e3 = PN(3, ('fell', 'down', '</s>')) # this is 0, so try less context
e2 = PN(2, ('down', '</s>')) # <= this works

# every time we tried a smaller N, we have to multiply by 0.4
# so we use a, then 0.4*0.4*b0, then 0.4*0.4*c1, then 0.4*d2, then 0.4*e2e.

print('a =', a, '\nb3=', b3, '\nb2=', b2, '\nb1=', b1, '\nb =', b)
print('c3=', c3, '\nc2=', c2, '\nc1=', c1, '\nd3=', d3, '\nd2=', d2, '\ne3=', e3, '\ne2=

print(a * (0.4*0.4*b) * (0.4*0.4*c1) * (0.4*d2) * (0.4*e2) )
print(P_stupid_backoff(3, ('<s>', 'The', 'grandstand', 'fell', 'down', '</s>')))

a = 0.005372992110137713
b3= 0
b2= 0
b1= 0
b = 8.611840246918683e-07
c3= 0
c2= 0
c1= 7.21757133431712e-05
d3= 0
d2= 0.01138101429453831
e3= 0
e2= 0.0011817757506744071

1.839836556015632e-20
1.839836556015632e-20

```

## Part C

Now we will implement the notion of *perplexity* as explained in lecture. Refer to the formula presented there to complete the following template, and verify that it passes all the tests.

```
In [68]: # Perplexity
def PP(N, W):
    probs = [P_stupid_backoff(N, W[max(i-N+1, 0):i+1]) for i in range(len(W))
    if 0 in probs:
        return float('inf')
    return np.prod([1/p for p in probs]) ** (1.0/(len(W)-1))
```

If we know that no probabilities can be 0, then P is same as P\_stupid\_backoff

```
In [69]: PP(2, ('<s>', 'The'))
```

```
Out[69]: 186.11603730316466
```

```
In [70]: P(2, ('<s>', 'The')) ** (-1/1)
```

```
Out[70]: 186.11603730316466
```

```
In [71]: P_stupid_backoff(2, ('<s>', 'The')) ** (-1/1)
```

```
Out[71]: 186.11603730316466
```

```
In [72]: PP(2, ('<s>', 'The', 'man', 'went'))
```

```
Out[72]: 308.91157551766736
```

```
In [73]: P(2, ('<s>', 'The', 'man', 'went')) ** (-1/3)
```

```
Out[73]: 308.91157551766736
```

```
In [74]: P_stupid_backoff(2, ('<s>', 'The', 'man', 'went')) ** (-1/3)
```

```
Out[74]: 308.91157551766736
```

```
In [75]: PP(2, ('<s>', 'The', 'man', 'went', 'to', 'the', 'house', '.', '</s>'))
```

```
Out[75]: 35.03849995731909
```

```
In [76]: P(2, ('<s>', 'The', 'man', 'went', 'to', 'the', 'house', '.', '</s>')) ** (-1/8)
```

```
Out[76]: 35.03849995731908
```

```
In [77]: P_stupid_backoff(2,('<s>','The','man','went','to','the','house','.', '</s>'))
```

```
Out[77]: 35.03849995731908
```

When probabilities may be 0, we must use stupid backoff

```
In [78]: PP(2,('<s>','where','is','the','grandstand'))
```

```
Out[78]: 7692.8065013245405
```

```
In [79]: P_stupid_backoff(2,('<s>','where','is','the','grandstand'))**(-1/4)
```

```
Out[79]: 7692.8065013245405
```

```
In [80]: PP(3,('<s>','where','is','the','grandstand'))
```

```
Out[80]: 8082.112259400885
```

```
In [81]: P_stupid_backoff(3,('<s>','where','is','the','grandstand'))**(-1/4)
```

```
Out[81]: 8082.112259400884
```

```
In [82]: PP(3,('<s>','The','grandstand','fell','down','</s>'))
```

```
Out[82]: 8852.055970670379
```

```
In [83]: P_stupid_backoff(3,('<s>','The','grandstand','fell','down','</s>'))**(-1/5)
```

```
Out[83]: 8852.055970670379
```

```
In [84]: PP(4,('<s>','The','grandstand','fell','down','</s>'))
```

```
Out[84]: 15339.392368477244
```

```
In [85]: P_stupid_backoff(4,('<s>','The','grandstand','fell','down','</s>'))**(-1/5)
```

```
Out[85]: 15339.392368477242
```

## Part D

Print out the first ten sentences in the training set, with their perplexities to 2 decimal places, using trigrams. Then do the same for the testing set.

Print out the text of the sentences in a readable form, e.g., for a sentence `w`, print it out using

```
' '.join(w[1:-1])
```

Notice the perplexities of the training set are generally smaller than the testing set!

```
In [86]: # your code here
# Define a function to print the sentence and its perplexity
def print_perplexity_and_sentence(dataset, N, count=10):
    for sentence in dataset[:count]:
        print(f"{PP(N, sentence):.2f} \t {' '.join(sentence[1:-1])}")

print_perplexity_and_sentence(train_sentences, 3, 10)
print("\n")
```

```
10.76    Muscle weakness did not improve , and the patient needed first a ca
ne , then crutches .
10.74    He replaced the flashlight where it had been stowed , got into his
own car and backed it out of the garage .
8.42     When he had given the call a few moments thought , he went into the
kitchen to ask Mrs. Yamata to prepare tea and sushi for the visitors , using
the formal English china and the silver tea service which had been donated t
o the mission , then he went outside to inspect the grounds .
10.04    -- On the basis of a differentiability assumption in function space
, it is possible to prove that , for materials having the property that the
stress is given by a functional of the history of the deformation gradients
, the classical theory of infinitesimal viscoelasticity is valid when the de
formation has been infinitesimal for all times in the past .
4.88     She said sharks have no bones and shrimp swam backward .
9.62     T. V. Barker , who developed the classification-angle system , was
about to begin the systematic compilation of the index when he died in 1931
.
10.32    He was then in man's hands .
32.57    4 .
15.57    `` Fifteen minutes , then ! !
11.27    Thus the cocktail party would appear to be the ideal system , but t
here is one weakness .
```

```
In [87]: # your code here
print_perplexity_and_sentence(test_sentences, 3, 10)
```

```
176.79 It is at least as important as the more dramatic attempts to break
down barriers of inequality in the South .
1306.74 the car's far windshield panel turned into a silver web wit
h a dark hole in the center .
69.22 `` I was just thinking how things have changed .
1071.75 She smiled , and the teeth gleamed in her beautifully model
ed olive face .
438.24 `` There isn't a chance of Myra's letting anything like that happen
.
174.63 On the other hand , many a pastor is so absorbed in ministering to
the intimate , personal needs of individuals in his congregation that he doe
s little or nothing to lead them into a sense of social responsibility and w
orld mission .
923.70 We live down by the Base commissary .
484.32 For example , the BBB has reported it was receiving four times as m
any inquiries about quack devices and 10 times as many complaints compared w
ith two years ago .
204.82 As a result , life had become a kind of continuous make-ready .
210.44 Some of the poems express a mood of joy in a newly discovered love
; ;
```

## Part E

Finally, we will find the perplexities of the the testing set with bigrams, trigrams, and quadrigrams. Complete the following template to verify that your results are consistent with the test results.

```
In [88]: # Find all the probabilities of the sentences in the testing set, multiply t
# and take the  $K^{th}$  root, where K is the number of tokens, excluding the
# So,  $K = (\text{sum of length of sentences}) - (\# \text{ of sentences})$ 

# We need to take the product of many small probabilities, so use math.log a
# underflow.

# Print out the perplexity as an integer.

import math

def compute_perplexity(N, sentences):
    log_prob = sum(math.log(P_stupid_backoff(N, s)) for s in sentences if P
    total_words = sum(len(s) - 1 for s in sentences)
    return int(math.exp(-log_prob / total_words))

for N in [2, 3, 4]:
    print(f"The perplexity of the testing set for {N}-grams is {compute_perp
```

The perplexity of the testing set for 2-grams is 387.  
 The perplexity of the testing set for 3-grams is 491.  
 The perplexity of the testing set for 4-grams is 935.

## Problem 3: Generative N-Gram Model

Now we will consider how to generate sentences using our N-gram model.

The idea is fairly simple. Suppose we have model using  $N=4$  (quadrigrams -- the algorithm for bigrams and trigrams is analogous):

1. To get  $w_1$ , choose a bigram (" $\langle s \rangle$ ",  $w_1$ ) randomly according the probability distribution stored in

$$\text{Ngram\_distribution}[2][ (" $\langle s \rangle$ ",  $w_1$ ) ].$$

1. To get  $w_2$ , choose a bigram (" $\langle s \rangle$ ",  $w_1, w_2$ ) randomly according the probability distribution calculated as

$$\text{Ngram\_distribution}[3][ (" $\langle s \rangle$ ",  $w_1, w_2$ ) ] / \text{Ngram\_distribution}[2][ (" $\langle s \rangle$ ",  $w_1$ ) ].$$

1. To get  $w_3$ , choose a trigram (" $\langle s \rangle$ ",  $w_1, w_2, w_3$ ) randomly according the probability distribution calculated as

$$\text{Ngram\_distribution}[4][ (" $\langle s \rangle$ ",  $w_1, w_2, w_3$ ) ] / \text{Ngram\_distribution}[3][ (" $\langle s \rangle$ ",  $w_1,$$$

1. Thereafter, for a sequence (" $\langle s \rangle$ ",  $w_1, w_2, \dots, w_{i-2}, w_{i-1}$ ), to get  $w_i$ , choose a quadrigram ( $w_{i-3}, w_{i-2}, w_{i-1}, w_i$ ) randomly according the probability distribution stored in

$$\text{Ngram\_distribution}[4][ (w_{i-3}, w_{i-2}, w_{i-1}, w_i) ] / \text{Ngram\_distribution}[3][ (w_{i-3}, w_{i-2},$$

1. When we generate the end of sentence marker `<\s>` we stop.

The problem is that this is difficult if we simply use the formulae given above: what we need is a separate probability distribution for each prefix.

## Part A

The first step, under the assumption that we are working with N-grams for  $N = 1, 2, 3$ , or 4, is to build a data structure that can sample from the distribution of next tokens given an  $(N-1)$ -gram of left context.

The best choice here is a nested default dictionary for N-grams for  $N = 2, 3, 4$ , the outer dictionary containing keys consisting of the first  $N-1$  tokens (we'll call this the *prefix*), with the value being an inner dictionary holding a probability distribution for the last token (we'll call this *wn*).

For this problem, you need to redo the construction of the list of N-grams and the distributions for each N, using the *entire* set of sentences, not just the testing set you used for the previous problems. You can easily do this by copying and pasting code from above.

```
In [89]: from collections import defaultdict, Counter
from nltk.corpus import brown
from nltk import ngrams

All_Ngrams = [None]
All_Ngram_distribution = [None]

for N in range(1, 5):

    ngram_list = [ngram for sent in brown.sents() for ngram in ngrams(sent,
All_Ngrams.append(ngram_list)

    distribution = defaultdict(Counter)
    for ngram in ngram_list:
        prefix, wn = tuple(ngram[:-1]), ngram[-1]
        distribution[prefix][wn] += 1
    for prefix, wn_counts in distribution.items():
        total_count = sum(wn_counts.values())
        distribution[prefix] = {wn: count / total_count for wn, count in wn_

    All_Ngram_distribution.append(distribution)
```



## Part B

Now we must build a data structure to solve the following problem: If we are working in an N-gram model for  $N = 2, 3$ , or  $4$ , given a sequence of tokens

$$\langle s \rangle \quad w_1 \ w_2 \ w_3 \ \cdots \ w_i$$

generate a sample word  $w_{i+1}$  using the distribution of the N-grams of the form

$$w_{i-N+1} \ \cdots \ w_{i-1} \ w_i \ w_{i+1}.$$

In other words, we use the last  $N - 1$  tokens of the sequence to determine a likely next token, given the distribution of N-grams starting with those  $N - 1$  tokens.

The best way to do this is to build a nested dictionary. Let us call the first  $N - 1$  tokens in an N-gram the *prefix* and the last token  $w_n$ . Then the outer dictionary is a `defaultdict` whose keys are the prefixes and values are an inner `defaultdict` whose keys are the  $w_n$  and whose values form a probability distribution for the ways that the prefix can be completed with a token  $w_n$ .

To give a simple example, suppose that in our corpus there are only the following bigrams whose first token is 'the':

```
('the', 'boy'),      ('the', 'baby'),      ('the', 'baby'),
('the', 'man')
```

Then our outer dictionary would have the prefix

```
('the',)
```

as a key, and the inner dictionary would store the probability that each of 'boy', 'baby', and 'man' would follow 'the', as shown in the next code cell.

Note that the prefix is an N-gram (a tuple) and  $w_n$  is simply a token.

```
In [90]: D = defaultdict(lambda: None)

D[('the',)] = defaultdict(lambda: 0)
D[('the',)][ 'boy' ] = 0.25
D[('the',)][ 'man' ] = 0.25
D[('the',)][ 'baby' ] = 0.5

D
```

```
Out[90]: defaultdict(<function __main__.<lambda>()>,
                    {('the',): defaultdict(<function __main__.<lambda>()>,
                                           {'boy': 0.25, 'man': 0.25, 'baby': 0.5})}})
```

Your task is to complete the following template and build a nested default dictionary giving the probability distributions for completions for (N-1)-grams.

Hint: For each N, you must create a `defaultdict` for all N-grams, whose key is the prefix and whose values are an inner `defaultdict`. For each N-gram, you should store the probability of the N-gram prefix+ $w_n$  under the key  $w_n$ . Then you must normalize these probabilities so that their sum is 1.0.

The end result will be that if you look up a prefix (N-1 tokens), you will have a probability distribution of possible  $w_n$  which can be used for the next token.

```
In [91]: from collections import defaultdict

def get_Ngram_dict(N):
    prefix_distribution = defaultdict(Counter)

    for ngram in All_Ngrams[N]:
        prefix, last_token = ngram[:-1], ngram[-1]
        prefix_distribution[prefix][last_token] += 1
    for prefix, wn_dict in prefix_distribution.items():
        total_count = sum(wn_dict.values())
        prefix_distribution[prefix] = {wn: count / total_count for wn, count in wn_dict.items()}

    return prefix_distribution

Ngram_nested_dict = [None]*5

for n in range(2,5):
    Ngram_nested_dict[n] = get_Ngram_dict(n)
```

```
In [92]: # tests: these should sum to (close to) 1.0

sum(Ngram_nested_dict[2][('<s>',)].values())
```

```
Out[92]: 1.00000000000000744
```

```
In [93]: sum(Ngram_nested_dict[3][('<s>', 'The')].values())
```

```
Out[93]: 0.99999999999999502
```

```
In [94]: sum(Ngram_nested_dict[4][('<s>', 'When', 'the')].values())
```

```
Out[94]: 1.00000000000000007
```

## Part C

Now that we have a way of sampling the next likely word, we will write a function which will predict the next word. You must sample from the probability distribution given a prefix, to choose a likely next word.

Hint: read about `numpy.random.choice`, in particular how you can set the parameter `p` to determine the probability of selecting a given key from the dictionary.

```
In [95]: def next_word(prefix):
          N = len(prefix) + 1

          if not (2 <= N <= 4) or not Ngram_nested_dict[N]:
              return None
          distribution_dict = Ngram_nested_dict[N].get(prefix, {})
          tokens, probabilities = list(distribution_dict.keys()), list(distribution_dict.values())

          return np.random.choice(tokens, p=probabilities) if probabilities else None
```

```
In [96]: # tests

          next_word((<s>,))
```

```
Out[96]: 'It'
```

```
In [97]: next_word((<s>, 'The'))
```

```
Out[97]: 'subjects'
```

```
In [98]: next_word((<s>, 'The', 'man'))
```

```
Out[98]: 'stood'
```

## Part D

Complete the following template to generate a random sentence by starting with the unigram `(<s>,)` and extending it by sampling until you generate the token `</s>`.

```
In [99]: # N is the parameter in N-gram

def generate_sentence(N):
    sentence = ['<s>']

    while sentence[-1] != '</s>':
        prefix = tuple(sentence[-(N-1):]) if len(sentence) >= N - 1 else tuple(sentence)
        sentence.append(next_word(prefix))

    return sentence
```

```
In [100]: # tests -- run this cell many times!

w1 = generate_sentence(2)

print(np.around(P(2,w1),2), ' '.join(w1[1:-1]),'\n')

w2 = generate_sentence(3)

print(np.around(P(3,w2),2), ' '.join(w2[1:-1]),'\n')

w3 = generate_sentence(4)

print(np.around(P(4,w3),2), ' '.join(w3[1:-1]),'\n')
```

63.1 Determine if we become so ' ', it , residence in the process .

15.18 He went into the pool too .

3.41 a middle distance often containing the major motif ; ;

## Part E

Experiment with generating sentences for various values of  $N = 2, 3, 4$ . How do the perplexities compare? Do you see a difference in the quality of the sentences? How well does it do with punctuation and quotes?

The typical view is that for larger values of  $N$ , the model is just "memorizing" the corpus. Do you think this is true? (You might look through the corpus to see what relationship your generated sentences, say for  $N = 4$ , have with the sentences in the corpus.

When i ran the test cell , i see that the quality of sentences are definitely changing . if the perplexity is high then it means that the randomness is high , meaning that the generated sentences are more inappropriate .

punctuations and quotes : for lower N gram (2), the punctuations & quotes are most random . but in the case of higher N gram (4), the punctuations and quotes are appropriate . but in 2-Gram also , at times(< 2%), the model is giving proper placement of punctuations .

i wrote a peice of code to check if the the generated sentences match the corpus . so i generated 100 sentences everytime to check this and more or less, i have to agree with the view of every one that "The typical view is that for larger values of N, the model is just "memorizing" the corpus." it is because , as said earlier , lower N gram model is producing based on assumptions and randomness where as the higher N gram is trying to match the sequence . Each time out of 100 for 4-Gram , there were an average of 18 sentences that matched the corpus while for the 2-Gram the average was 2

In [ ]:

In [ ]: