

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 \(https://www.youtube.com/channel/UCfSqNB0yh99yuG4p4nzjPOA\)](https://www.youtube.com/channel/UCfSqNB0yh99yuG4p4nzjPOA).

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.

chatGPT: requested chatGPT to debug `(prob_not_normalized ** -1/(N-1))` vs. `pow(prob_not_normalized, -1/(N-1))` in 2C and clarify how Python deals with floating point number computation.
Python documentation: `np.random.seed`, `np.random.choice`, `dict.get()`, `assert` or other error handling techniques

```
In [2]: 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 /Users/yfsong/nltk_data...
[nltk_data] Package brown is already up-to-date!
```

Out[2]: True

```
In [3]: brown.sents()
```

Out[3]: [['The', 'Fulton', 'County', 'Grand', 'Jury', 'said', 'Friday', 'an', 'investigation', 'of', 'Atlanta's', 'recent', 'primary', 'election', 'produced', '``', 'no', 'evidence', '``', 'that', 'any', 'irregularities', 'took', 'place', '.'], ['The', 'jury', 'further', 'said', 'in', 'term-end', 'presentments', 'that', 'the', 'City', 'Executive', 'Committee', ',', 'which', 'had', 'over-all', 'charge', 'of', 'the', 'election', ',', '``', 'deserves', 'the', 'praise', 'and', 'thanks', 'of', 'the', 'City', 'of', 'Atlanta', '``', 'for', 'the', 'manner', 'in', 'which', 'the', 'election', 'was', 'conducted', '.'], ...]

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).

3. Then split `shuffled_sentences` into sets `training_sents` (first 99.9% of the sentences) and `testing_sents` (last 0.1%).
4. Print out the length of the training and testing sets.
5. 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 [4]: # First, shuffle the set of sentences
import copy

seed(0) # set the seed to ensure reproducibility

# your code here
sentences = list(brown.sents()) # original list
shuffled_sentences = sentences.copy()
shuffle(shuffled_sentences) # randomized list: np.random.shuffle does not return any value so no assignment ne
split_at = math.floor(len(shuffled_sentences)*0.999)
train, test = shuffled_sentences[:split_at], shuffled_sentences[split_at:]
```

```
In [5]: print(f"Length training set: {len(train)}.")
print(f"Length testing set: {len(test)}.\n")
print(f"Start of training set:\n{train[0]}\n")
print(f"Start of testing set:\n{test[0]}")
```

Length training set: 57282.
Length 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', 'd
own', '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.

In [6]: *# put ``<s>`` at beginning and ``</s>`` at end of all sentences.*

```
def bracket_sentence(sent):
    return ['<s>'] + sent + ['</s>']

# your code here
train_markers_added = [bracket_sentence(s) for s in train]
test_markers_added = [bracket_sentence(s) for s in test]
print(f"Start of training set:\n{train_markers_added[0]}\n")
print(f"Start of testing set:\n{test_markers_added[0]}\n")
```

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 [7]: *# 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), etc.

```
def get_Ngrams_for_sentence(N,s): # N for N-grams, s for one sentence
    assert 1<= N <= 4 and isinstance(N, int)

    return [tuple(s[i:i + N]) for i in range(len(s) - N + 1)]
```

your code here

```
print(f"Unigrams (N=1):\n{get_Ngrams_for_sentence(1, train_markers_added[0])}\n")
print(f"Bigrams (N=2):\n{get_Ngrams_for_sentence(2, train_markers_added[0])}\n")
print(f"Trigrams (N=3):\n{get_Ngrams_for_sentence(3, train_markers_added[0])}\n")
print(f"Tetragrams (N=4):\n{get_Ngrams_for_sentence(4, train_markers_added[0])}\n")
```

Unigrams (N=1):

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

Bigrams (N=2):

```
[('<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>')]
```

Trigrams (N=3):

```
[('<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>')]
```

Tetragrams (N=4):

```
[('<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 [8]: Ngrams = [[] for _ in range(5)]
# first slot is empty, then Ngram[1] will hold unigrams, Ngram[2] will hold bigrams, etc.
# [None] + [[]] would make all n-grams ref to the same memory location

for i in range(1,5):
    for s in train_markers_added:
        Ngrams[i] += get_Ngrams_for_sentence(i,s)

num_ngrams = [None] + [len(Ngrams[i]) for i in range(1,5)]
first_5 = [None] + [Ngrams[i][:5] for i in range(1,5)]

for i in range(1,5):
    print(f"There are {num_ngrams[i]} N-grams in Ngrams[{i}] and the first 5 are:\n{first_5[i]}\n")
```

There are 1274667 N-grams in Ngrams[1] and the first 5 are:

```
[('<s>',), ('Muscle',), ('weakness',), ('did',), ('not',)]
```

There are 1217385 N-grams in Ngrams[2] and the first 5 are:

```
[('<s>', 'Muscle'), ('Muscle', 'weakness'), ('weakness', 'did'), ('did', 'not'), ('not', 'improve')]
```

There are 1160103 N-grams in Ngrams[3] and the first 5 are:

```
[('<s>', 'Muscle', 'weakness'), ('Muscle', 'weakness', 'did'), ('weakness', 'did', 'not'), ('did', 'not', 'improve'), ('not', 'improve', ',')]
```

There are 1102821 N-grams in Ngrams[4] and the first 5 are:

```
[('<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 [288]: # Create a defaultdict with the frequency distribution for the training set for a given N.
from collections import defaultdict

def get_Ngram_distribution(N,Ngrams): # N for 1-4, Ngrams for collections of unigrams, etc
    new_dict = defaultdict(lambda: 0)
    for t in Ngrams: # frequency distribution
        new_dict[t] += 1
    for k in new_dict:
        new_dict[k] /= num_ngrams[N]
    return new_dict

# probability distribution
Ngram_distribution = defaultdict(int)
for i in range(1,5):
    Ngram_distribution[i] = get_Ngram_distribution(i,Ngrams[i])

# Since all sentences will start with <s>, its probability should be 1.0
Ngram_distribution[1][('<s>',)] = 1.0
```

```
In [291]: # tests
for i in range(1,5):
    print(f"The total number of unique Ngrams[{i}] is {len(Ngram_distribution[i].keys())}.\n")

print(f"\nThe Probability of ('<s>',) is {np.around(Ngram_distribution[1][('<s>',)],8)}.")
print(f"\nThe Probability of ('to',) is {np.around(Ngram_distribution[1][('to',)],8)}.")
print(f"\nThe Probability of ('to','the') is {np.around(Ngram_distribution[2][('to','the')],8)}.")
print(f"\nThe Probability of ('to','the','house') is {np.around(Ngram_distribution[3][('to','the','house')],8)}")
print(f"\nThe Probability of ('to','the','house','.') is {np.around(Ngram_distribution[4][('to','the','house',
```

The total number of unique Ngrams[1] is 56029.

The total number of unique Ngrams[2] is 454055.

The total number of unique Ngrams[3] is 877117.

The total number of unique Ngrams[4] is 1029391.

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

and explicated in this StackOverflow post:

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

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

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

Finally, let

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

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 definiti
on

```

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

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 [292]: *# Probability of a list of tokens using N-grams*

```
# Calculate the probability of the last token in W, as we did in the calculation of p1, p2, etc. above
# check numerator: if it is 0, return 0 immediately and do not do the division (to avoid possible
# division by 0)
```

```
# N == len(W)
```

```
def PN(N, W):
    # W: a list of tokens; N == len(W)
    assert N>0, "N must be positive."

    if N == 1: # only one token
        return Ngram_distribution[1].get((W[0],), 0)

    numer = Ngram_distribution[N].get(tuple(W), 0)
    denom = Ngram_distribution[N-1].get(tuple(W[:-1]), 0)

    if numer != 0 and denom != 0:
        return numer/denom
    return 0
```

```
# Now calculate for a whole sequence: use PN(..) for bigram, trigram, etc. up to size of
# model, then slide PN(...) across the sequence, as shown above.
```

```
def P(N, W):
    assert N<=len(W), "N cannot be greater than length of W."

    base = 1
    # Pr from 1-gram to (N-1)-gram: p1, ..., p(N-1)
    for n in range(1, N):
        each_prob = PN(n, W[:n])
        if each_prob == 0:
            return 0
        base *= each_prob

    # N grams for slices of tokens of size N
    for i in range(len(W)-N+1):
        each_tuple = W[i:i+N]
        each_prob = PN(N, each_tuple)
        if each_prob == 0:
            return 0
        base *= each_prob

    return base
```

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

In [293]: *# 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 [294]: # 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 [295]: # 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 [296]: # Sentence: I love NLP
# using trigrams

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

a = Ngram_distribution[2][('<s>', 'I')]
b = Ngram_distribution[3][('<s>', 'I', 'love')] / Ngram_distribution[2][('<s>', 'I',)]
c = Ngram_distribution[3][('I', 'love', 'NLP')] / Ngram_distribution[2][('I', 'love')]
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*d*d1)
print('P(3,...): ', P(3, ('<s>', 'I', 'love', 'NLP', '</s>'))))

a= 0.001128648701930778
b= 0.0015274769289326804
c= 0.0
d1= 0
d2= 0

a*b*d*d1: 0.0
P(3,...): 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 [343]: # Probability with stupid backoff
# same as previous, but have to use recursive (or iterative) method instead of
# 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-2, etc. down to 1 token.

# Assumes W is a tuple

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

def PN_with_stupid_backoff(N, W):
    # N == len(W); W: a single N-gram
    # a recursive method that returns backed-off probability for single N-gram

    prob = PN(N, W)

    if prob > 0:
        return prob

    if N == 1:
        return brown.words().count(W[0]) / len(brown.words())

    return 0.4 * PN_with_stupid_backoff(N-1, W[1:])

# Note that for training set, this will be same as P(N,W) since all probabilities are non-zero

# Really the same as P(N,W) except using the previous function, and no need to check for 0,
# since it never returns 0 for a sequence using words from the corpus

def P_stupid_backoff(N,W):
    # assert N<=len(W), "N cannot be greater than length of W."

    base = 1
    # Pr from 1-gram to (N-1)-gram: p1, ..., p(N-1)
    for n in range(2, N):
        each_prob = PN_with_stupid_backoff(n, W[:n])
        base *= each_prob

    # N grams for slices of tokens of size N
    for i in range(len(W)-N+1):
        each_tuple = W[i:i+N]
        each_prob = PN_with_stupid_backoff(N, each_tuple)
        base *= each_prob

    # print(f"final probability is {base}")
    return base
```

```
In [ ]: # saved
def PN_with_stupid_backoff(N, W):
    # N == len(W); W: a single N-gram
    # a recursive method that returns backed-off probability for single N-gram
    assert N == len(W), "N must be equal to length of W"

    numer = Ngram_distribution[N].get(tuple(W), 0)
    denom = Ngram_distribution[N-1].get(tuple(W[:-1]), 0) if N > 1 else 1

    if numer and denom:
        return numer/denom

    if N == 1:
        return Ngram_distribution[1].get(W, 1/num_all_tokens)
    return 0.4 * PN_with_stupid_backoff(N-1, tuple(W[1:]))
```

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

```
In [298]: # 'grandstand' is in testing set but not in the training set
# This uses bigrams

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

Out[298]: 0

```
In [299]: 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 [300]: # 'grandstand' is in testing set but not in the training set
# This uses trigrams
```

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

```
Out[300]: 0
```

```
In [301]: 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 whole corpus
d = list(brown.words()).count('grandstand') / len(brown.words())      # <= this works

print('a =',a,'\nb3=',b3,'\nb2=',b2,'\nc =',c,'\nd3=',d3,'\nd2=',d2,'\nd1=',d1,'\nd =',d,'\n')

# 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= 2.3535558698860173e-06
d = 8.611840246918683e-07
```

```
2.3436917228933544e-16
2.3436917228933544e-16
```

```
In [302]: a = PN(2,('<s>','The')) # <= this works

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

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=',e2,'\n')

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 [325]: # Perplexity

# We assume that W starts with <s>, may not end with </s>

def PP(N,W):
    prob_not_normalized = P_stupid_backoff(N, W)
    print(f"unnormalized pr is {prob_not_normalized}")

    N = len(W)

    perplexity = pow(prob_not_normalized, -1/(N-1)) # assume that W starts with <s>
    '''
    used chatGPT here to help me debug
    (prob_not_normalized ** -1/(N-1)) vs. pow(prob_not_normalized, -1/(N-1))
    in regards of floating point number handling of python
    and also: https://stackoverflow.com/questions/20969773/exponentials-in-python-xy-vs-math-powx-y
    '''

    return perplexity
```

If we know that no probabilities can be 0, then P is same as P_stupid_backoff

```
In [304]: PP(2,('<s>','The'))

unnormalized pr is 0.005372992110137713
```

```
Out[304]: 186.11603730316466
```

```
In [305]: P(2,('<s>','The'))**(-1/1)
```

```
Out[305]: 186.11603730316466
```

```
In [306]: P_stupid_backoff(2,('<s>','The'))**(-1/1)
```

```
Out[306]: 186.11603730316466
```

```
In [307]: PP(2, ('<s>', 'The', 'man', 'went'))
          unnormalized pr is 3.3923249973453404e-08
Out[307]: 308.91157551766736

In [308]: P(2, ('<s>', 'The', 'man', 'went'))**(-1/3)
Out[308]: 308.91157551766736

In [309]: P_stupid_backoff(2, ('<s>', 'The', 'man', 'went'))**(-1/3)
Out[309]: 308.91157551766736

In [310]: PP(2, ('<s>', 'The', 'man', 'went', 'to', 'the', 'house', '.', '</s>'))
          unnormalized pr is 4.401857289920249e-13
Out[310]: 35.03849995731908

In [311]: P(2, ('<s>', 'The', 'man', 'went', 'to', 'the', 'house', '.', '</s>'))**(-1/8)
Out[311]: 35.03849995731908

In [312]: P_stupid_backoff(2, ('<s>', 'The', 'man', 'went', 'to', 'the', 'house', '.', '</s>'))**(-1/8)
Out[312]: 35.03849995731908

When probabilities may be 0, we must use stupid backoff

In [313]: PP(2, ('<s>', 'where', 'is', 'the', 'grandstand'))
          unnormalized pr is 2.855359302895301e-16
Out[313]: 7692.8065013245405

In [314]: P_stupid_backoff(2, ('<s>', 'where', 'is', 'the', 'grandstand'))**(-1/4)
Out[314]: 7692.8065013245405

In [315]: PP(3, ('<s>', 'where', 'is', 'the', 'grandstand'))
          unnormalized pr is 2.3436917228933544e-16
Out[315]: 8082.112259400884

In [316]: P_stupid_backoff(3, ('<s>', 'where', 'is', 'the', 'grandstand'))**(-1/4)
Out[316]: 8082.112259400884

In [317]: PP(3, ('<s>', 'The', 'grandstand', 'fell', 'down', '</s>'))
          unnormalized pr is 1.839836556015632e-20
Out[317]: 8852.055970670379

In [318]: P_stupid_backoff(3, ('<s>', 'The', 'grandstand', 'fell', 'down', '</s>'))**(-1/5)
Out[318]: 8852.055970670379

In [319]: PP(4, ('<s>', 'The', 'grandstand', 'fell', 'down', '</s>'))
          unnormalized pr is 1.1774953958500048e-21
Out[319]: 15339.392368477242

In [320]: P_stupid_backoff(4, ('<s>', 'The', 'grandstand', 'fell', 'down', '</s>'))**(-1/5)
Out[320]: 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 [321]: *# first 10 sents in training data, .2f perplexities, N=3*

```
import re
for i in range(10):
    sent = train_markers_added[i]
    perplex = PP(3, sent)
    sent_to_print = ' '.join(sent[1:-1])
    print(f"[{i+1}] {sent_to_print}\nPerplexity: {perplex:.2f}\n")
```

unnormalized pr is 2.6776408392894215e-19

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

Perplexity: 10.76

unnormalized pr is 1.7995728627929014e-25

[2] He replaced the flashlight where it had been stowed , got into his own car and backed it out of the garage .

Perplexity: 10.74

unnormalized pr is 1.8174843330080426e-53

[3] 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 to the mission , then he went outside to inspect the grounds .

Perplexity: 8.42

unnormalized pr is 7.576209061206683e-63

[4] -- 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 deformation has been infinitesimal for all times in the past .

Perplexity: 10.04

unnormalized pr is 5.5388151774745665e-09

[5] She said sharks have no bones and shrimp swam backward .

Perplexity: 4.88

unnormalized pr is 2.817712069632965e-27

[6] T. V. Barker , who developed the classification-angle system , was about to begin the systematic compilation of the index when he died in 1931 .

Perplexity: 9.62

unnormalized pr is 7.793525348269414e-09

[7] He was then in man's hands .

Perplexity: 10.32

unnormalized pr is 2.8945751201105465e-05

[8] 4 .

Perplexity: 32.57

unnormalized pr is 2.9017909506161495e-10

[9] `` Fifteen minutes , then ! !

Perplexity: 15.57

unnormalized pr is 1.0382352039112117e-20

[10] Thus the cocktail party would appear to be the ideal system , but there is one weakness .

Perplexity: 11.27

```
In [334]: # first 10 sents in testing data, .2f perplexities, N=3
```

```
for i in range(10):
    sent = test[i]
    perplex = PP(3, sent)
    sent_to_print = ' '.join(sent)
    print(f"[{i+1}] {sent_to_print}\nPerplexity: {perplex:.2f}\n")
```

unnormalized pr is 2.1514445098661873e-47

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

Perplexity: 215.46

unnormalized pr is 5.763932045997189e-56

[2] the car's far windshield panel turned into a silver web with a dark hole in the center .

Perplexity: 1775.70

unnormalized pr is 1.0242617660756136e-18

[3] `` I was just thinking how things have changed .

Perplexity: 99.73

unnormalized pr is 4.513785190252695e-43

[4] She smiled , and the teeth gleamed in her beautifully modeled olive face .

Perplexity: 1808.60

unnormalized pr is 2.0928294474577326e-35

[5] `` There isn't a chance of Myra's letting anything like that happen .

Perplexity: 776.14

unnormalized pr is 8.049929993566162e-96

[6] On the other hand , many a pastor is so absorbed in ministering to the intimate , personal needs of individuals in his congregation that he does little or nothing to lead them into a sense of social responsibility and world mission .

Perplexity: 183.72

unnormalized pr is 8.751588919085126e-24

[7] We live down by the Base commissary .

Perplexity: 1967.83

unnormalized pr is 1.0315295626773937e-80

[8] For example , the BBB has reported it was receiving four times as many inquiries about quack devices and 10 times as many complaints compared with two years ago .

Perplexity: 573.00

unnormalized pr is 2.4449106550837164e-29

[9] As a result , life had become a kind of continuous make-ready .

Perplexity: 242.28

unnormalized pr is 1.8083086767545365e-36

[10] Some of the poems express a mood of joy in a newly discovered love ; ;

Perplexity: 241.46

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 [340]: # Find all the probabilities of the sentences in the testing set, multiply them,
# and take the  $K^{\text{th}}$  root, where  $K$  is the number of tokens, excluding the <s> tokens.
# 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 and math.exp to avoid
# underflow.

# Print out the perplexity as an integer.

import math

def all_perplexity(N, W):

    total_log_prob = 0
    num_wrds = sum([len(w) for w in W])
    K = num_wrds - len(W)
    # total number of words in all sentences, but subtracts one word <s> for each sentence

    for w in W:
        each_prob = P_stupid_backoff(N, w)

        if each_prob <= 0:
            each_prob = 1e-10 # make sure no log(0)

        total_log_prob += math.log(each_prob)

    perplexity = math.exp(-total_log_prob/K)

    return perplexity
```

```
In [ ]: for i in range(2,5):
        print(f"The perplexity of the testing set for {i}-grams is {all_perplexity(i, test)}")
```

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 (" <s> ", w_1) randomly according the probability distribution stored in

$$\text{Ngram_distribution}[2][(" <s> ", w_1)].$$

2. To get w_2 , choose a bigram (" <s> ", w_1, w_2) randomly according the probability distribution calculated as

$$\text{Ngram_distribution}[3][(" <s> ", w_1, w_2)] / \text{Ngram_distribution}[2][(" <s> ", w_1)].$$

3. To get w_3 , choose a trigram (" <s> ", w_1, w_2, w_3) randomly according the probability distribution calculated as

$$\text{Ngram_distribution}[4][(" <s> ", w_1, w_2, w_3)] / \text{Ngram_distribution}[3][(" <s> ", w_1, w_2)].$$

4. Thereafter, for a sequence (" <s> ", $w_1, w_2, \dots, w_{i-2}, w_{i-1}$), 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

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

5. 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 [113]: from collections import defaultdict

All_Ngrams = [None] + [[] for _ in range(4)]

shuffled_sents_markers = ['<s>'] + s + ['</s>'] for s in shuffled_sentences]

for i in range(1,5):
    for s in shuffled_sents_markers: # shuffled_sents_markers: set of all sentences
        All_Ngrams[i] += get_Ngrams_for_sentence(i, s)

all_num_ngrams = [None] + [len(All_Ngrams[i]) for i in range(1,5)]

All_Ngram_distribution = [None]

for i in range(1,5):
    dtb = defaultdict(int)
    for ngram in All_Ngrams[i]: # frequency distribution
        dtb[ngram] += 1
    for k in dtb: # normalize freq to get prob
        dtb[k] /= all_num_ngrams[i]
    All_Ngram_distribution.append(dtb)

All_Ngram_distribution[1][('<s>',)] = 1.0
```

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 [92]: 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[92]: 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 [149]: # now build a dictionary of lists of completions, with probabilities
# For N-gram, key is prefix, of length N-1, which returns a dictionary
# with keys that are tokens (the last token wn in the N-gram), with
# values the number of times that N-gram (prefix),wn occurs.

# N here is the length of the whole N-gram, so the prefix is of length N-1

def get_Ngram_dict(N):
    prefix_dict = defaultdict(lambda: defaultdict(float))

    for ngram, prob in All_Ngram_distribution[N].items():
        prefix = ngram[:-1]
        wn = ngram[-1]
        prefix_dict[prefix][wn] = prob

    for prefix, next_dict in prefix_dict.items():
        sum_prob = sum(next_dict.values())
        for w, prob in next_dict.items():
            prefix_dict[prefix][w] = prob / sum_prob
            # normalize to ensure the sum of probs for all potential next words given the prefix is 1

    return prefix_dict

Ngram_nested_dict = [None]*5

for n in range(1,5): # Also include unigrams for completeness
    Ngram_nested_dict[n] = get_Ngram_dict(n)
```

```
In [150]: # tests: these should sum to (close to) 1.0
```

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

```
Out[150]: 1.00000000000002396
```

```
In [151]: sum(Ngram_nested_dict[3][('<s>', 'The')].values())
```

```
Out[151]: 0.9999999999999709
```

```
In [152]: sum(Ngram_nested_dict[4][('<s>', 'When', 'the')].values())
```

```
Out[152]: 1.0000000000000013
```

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 [153]: # given a prefix, randomly choose next token using the appropriate probability distribution
```

```
# len(prefix) must be 1, 2, 3, or 4
```

```
def next_word(prefix):
    N = len(prefix) + 1
    # given a prefix of length x,
    # we want to get the distribution of n-grams with token of length x+1

    assert prefix in Ngram_nested_dict[N], "prefix must be in the data"

    next_words = list(Ngram_nested_dict[N][prefix].keys()) # list of next probable word
    probs = list(Ngram_nested_dict[N][prefix].values()) # list of corresponding probs

    probs_normalized = [prob / sum(probs) for prob in probs]

    word_to_choose = np.random.choice(next_words, p=probs_normalized)

    return word_to_choose
```

```
In [172]: # tests
```

```
next_word('<s>')
```

```
Out[172]: 'The'
```

```
In [170]: next_word('<s>', 'The')
```

```
Out[170]: 'student'
```

```
In [167]: next_word(('<s>', 'The', 'man'))
```

```
Out[167]: 'was'
```

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 [345]: # N is the parameter in N-gram

def generate_sentence(N):
    sent = ['<s>'] # starts with the unigram ('<s>',)

    while 1:
        prefix = tuple(sent[-(N-1):])
        word = next_word(prefix)
        sent.append(word)
        if word == '</s>':
            break

    return sent
```

```
In [352]: # tests -- run this cell many times!
w1 = generate_sentence(2)
print(np.around(PP(2,w1),2), ' '.join(w1[1:-1]),'\n\n')
w2 = generate_sentence(3)
print(np.around(PP(3,w2),2), ' '.join(w2[1:-1]),'\n\n')
w2 = generate_sentence(3)
print(np.around(PP(3,w2),2), ' '.join(w2[1:-1]),'\n\n')
```

unnormalized pr is 4.5147377960729765e-72

84.77 that God through small business with a little extra work at the commerce '' pas de l'Independance) Baldwin Longstreet and squash keep the Descent were organized whole sera , designed to straight to ! !

unnormalized pr is 9.425272445774055e-12

5.43 She thought royal status might come easy to go home by bus ? ?

unnormalized pr is 6.708535219504624e-40

6.82 Mike struck with the peripheral resolution losses resulting with a girl about ten days were at work she was three parts , it can claim this promise tendered by the National Defense , Charles Hitch , who was `` dead game sport '' in diameter .

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.

**

Regarding perplexities, it seems that the perplexity varies for different Ns, but a stable observation is that perplexity for bigrams is relatively higher compared to the ones generated using trigrams and quadragrams. Since a model with lower perplexity is considered "better", the results indicate that trigram models may have a slightly better performance than quadragram models and largely better than bigram ones.

As N increases, the quality of sentences appear to be higher, the position of punctuations are more reasonable, though I think for all Ns we tested the punctuations remain a gray area that we may need more optimizations: trigram and quadragram models produce more coherent and understandable sentences.

In theory, as N goes up, we give more contexts and the sentences generated probabilistically resemble the ones in the training data. Therefore, the model indeed memorizes the combination of tokens.**