

ADDIS ABABA UNIVERSITY

ADDIS ABABA INSTITUTE OF TECHNOLOGY

SCHOOL OF INFORMATION TECHNOLOGY AND ENGINEERING – SITE

Department of MSC in Artificial Intelligence

Course Assignment

Course: NLP

Name: Mintesnot Fikir

Section: Regular

IDs: GSR1669/15

Submitted to: Fantahun B. (PhD)

Submission date November 2023

NLP Assignment-1 (12%)

Fantahun Bogale • Oct 20 (Edited Nov 6) 12 points Due Yesterday, 11:59 PM

- This assignment has three (3) tasks each having its own weight. The total weight of this assignment will be 12%.
- This assignment is based on one version of the Amharic Corpus (GPAC) created by this paper. <https://doi.org/10.3390/info12010020>.
- You can access this corpus using this link: link.
[https://drive.google.com/file/d/1jEOAORyITBAiNKes_JerD6gFF7cnLFyr/view?
usp=sharing](https://drive.google.com/file/d/1jEOAORyITBAiNKes_JerD6gFF7cnLFyr/view?usp=sharing)
- This assignment is an individual assignment. Experience sharing and discussion is possible, however each student is responsible to submit and defend his work.
- Presentation may be required (get ready for that) >

Note: The corpus may require a good size of RAM to load; hence you can make use of Google Collab, Kaggle or more if your device is incapable. > Questions:

1 N-gram language model (6%)

- 1.1 Create n-grams for n=1, 2, 3, 4. You can show sample prints.
- 1.2 Calculate probabilities of n-grams and find the top 10 most likely n-grams for all n.
- 1.3 What is the probability of the sentence. "እ.ት.የአድ.የ.ታ.ቃ.ቅ.ሁ.በ.ኋ.ት". You can also try more sentences.
- 1.4 Generate random sentences using n-grams; explain what happens as n-increases based on your output. >

2 Evaluate these Language Models Using Intrinsic Evaluation Method (3%)

3 Evaluate these Language Models Using Extrinsic Evaluation Method (3%)

You can make use of any task convenient for you to evaluate the n-gram models you have created.

Solution

Outline

Preprocessing the data

1. N-gram language mode

- 1.1. Create n-grams for n=1, 2, 3, 4. You can show sample prints.
- 1.2. Calculate probabilities of n-grams and find the top 10 most likely n-grams for all n.
- 1.3. What is the probability of the sentence. "አትናኝ ታሪክዎች ስነ የት". You can also try more sentences.
- 1.4. Generate random sentences using n-grams; explain what happens as n-increases based on your output

2. Evaluate these Language Models Using Intrinsic Evaluation Method

3. Evaluate these Language Models Using Extrinsic Evaluation Method

Preprocessing

Import necessary libraries

In [1]:

```
import pickle
import nltk
import random
import gc
from nltk import ngrams
from nltk.util import ngrams
from nltk import FreqDist
from collections import Counter
import random
import math
import time
import string
```

Reading the corpus from local folder(text file) as a text file

```
In [2]: starttime = time.time()
with open('/kaggle/input/amharic-corpus-general/GPAC.txt', 'r', encoding='utf-8') as file:
    text = file.read()
endtime = time.time()
print(f"Time taken for read the text as a string is : {(endtime - starttime)/60}")
```

Time taken for read the text as a string is : 0.3450124700864156 minuets

```
In [3]: print(type(text))
```

As we have imported text file as a single string, next step is to make it comfortable for further processing, so we need to tokenize the string into list of tokens

```
In [4]: print("\n\n Sample text data:\n\n",text[4:3000])
```

Sample text data:

The dataset contains numerous extraneous characters and symbols which require further processing. Below is a list of some of these:

Tokenize the corpus

To facilitate more efficient processing of the corpus, the initial phase should involve tokenization. In this case, since the text comprises a mix of individual characters, symbols, and potential punctuation marks, tokenization would involve categorizing these elements into meaningful groups that can be separately identified and analyzed.

By breaking down the text into tokens, we can apply subsequent processing steps more effectively. These steps could include filtering out unnecessary symbols, converting characters to a uniform case for consistency, or identifying and preserving meaningful punctuation for further linguistic analysis. Tokenization simplifies these tasks by providing a structured approach to dissecting the text, allowing for a more manageable and targeted processing workflow.

we have 2 ways of tokenization either create Custom Tokenization Function from scratch or use third party library

```
In [5]: # tokenizer functions
import string
def tokenizer(text: str):
    tokens = []
    current_token = ''

    for char in text:
        if char.isspace() or char in string.punctuation:
            if current_token:
                tokens.append(current_token)
                current_token = ''
        else:
            current_token += char

    if current_token:
        tokens.append(current_token)

    return tokens
```

- The above custom tokenization function tokenizer iterates over each character in the input string. It builds tokens by appending characters until it encounters a whitespace or a punctuation character, at which point it adds the current token to the list and starts building the next one. This function does not use regular expressions and relies purely on sequential character checking, making it straightforward but potentially less efficient for large texts compared to regex-based solutions.

```
In [6]: start_time = time.time()
tokens = tokenizer(text)
```

```

end_time = time.time()
print(f"\n\nTime taken for tokenization: {(end_time - start_time)/60} minutes \n")
print(tokens[:200])
print('\nSize of valid tokens are: ',len(tokens), 'Words(tokens)')

```

Time taken for tokenization: 2.299083085854848 minutes

```

['ምን', 'መስተኞች', 'አጠቃላይ', 'አጥቃለሁ', 'ጥረዋ', 'እርዳት', 'ልተታደሙ', 'ያልቻሉት
ው', 'የእለም', 'የእግዥ', 'ኋስ', 'ዋ', 'ለ19ኛ', 'ኋክ', 'በደብብ', 'እፍሪ', 'ሰመጣ', 'በፌዴ', 'እየ
ቻ', 'እንደሸ', 'ባረሪ', 'ልክ', 'በአመቱ', 'በለስ', 'ቀጥና', 'ልል', 'ዋ', 'ልተታደም', 'ሁለተ', 'ልቻ
?', 'ወደ', 'ደብብ', 'አፍሪ', 'አከተ፡ይናው', 'በግ', 'በራዘዎ', 'አፍሪ', 'አብር', 'የመኖር', 'ወድድ
ር', 'በደብብ', 'አፍሪ', 'ተካሂወል፡', 'ከተለያ', '14', 'የእፍሪ', 'አገተ', 'የተመማበ', '26', 'የህ
ል', 'ተሳታፊወቻ', 'የተከፈለበት', 'ይህ', 'ወድድር', 'ግለሰቦች', 'በደታች', 'ሁኔታ', 'ወሰጥ', 'በማለፍ',
'በቻታቻውን', 'የሚያስመስክሱት', 'መሆኑ', 'ልምተናል፡', 'የሚገጥማቻውን', 'የተለያ', 'ኤተናወቻ', 'በትላቅ
ትና', 'በጥበብ', 'ማለፍ፡', 'ከለሞቻ', 'ታርጉቶች', 'መዘላቹ፡', 'ቃጉቶች', 'በበለጠት', 'መፍታት',
'ወዘተ', 'በየነዚው', 'ከሚደረገው', 'ቀነሳ', 'ተርጉሙ', 'ለ91', 'ዋናት', 'የህል', 'በወድድኑ', 'መቆየት',
'የቻል', 'ሁለት', 'ተመክሮቻዊቻ', 'አያንዳንዶቻው', '200', 'ይህ', 'ዳለር', 'አንድማስለመው', 'በኅንር', 'ነብ
ር፡', 'በዘመርመው', 'ወድድር', 'አገተቻ', 'ኋኙ', 'አና', 'የትባለ', 'ሁለት', 'ወጥቶቻ', 'በቻስል
ዋ', 'ኋኙ', 'ቀደም', 'በለስ', 'የቅናወ', 'በለስ', 'ቢሆን', 'የኋይም', 'በቅርቡ', 'ከወድድር', 'ወመብ', 'ሁና
ለቻ፡ይህቻ', 'የአገተ', 'በጥቻ', 'ተስፋ', 'ወደ', 'አስናወለት', 'ለማሽንር', 'የህዝብ', 'የድም', 'የ
መናሽ', 'መሆኑ', 'የተገኘበው', 'ወይም', 'ኋው', 'አንዳሸ', '835', 'የሚል', 'አገዋ', 'ተሬ', 'የስተላለ
ልወቻ', 'ይህ', 'የኋኙ', 'ከወድድኑ', 'ከመሰናበት', 'በፈት፡ወይም', 'የአገተ', 'የህል', 'በኢናዲነት', 'የሚሰጠሩ
ት', 'ከፍተቻ', 'ገኘት፡', 'አገራ', 'አንዳተሻነ', 'የሚል', 'ከፍተቻ', 'በጋት', 'አንዳረሰበት', 'ይሰማቻል፡',
'ገኘት', 'በይሁን', 'በጋት', 'የወይም', 'የዋህነት', '፡፡የኋይም', 'አጥቃለያም', 'ይሰነኩ', 'ይሁን', 'በሚል',
'አንዳሸ', 'ከገተ', 'በጋት', 'የሚያገበባቸውን', 'አገር', 'ወይም', 'ዘጋጀ', 'ኋኙ', 'የዋህቻ', 'አለቻዋለሁ፡የዋ
ሁቻ', 'ሁሬ', 'አትስት', 'በለስ', 'የኋይም', 'በለስ', 'አጥቃለያም', 'አትስት፡', 'ወድድኑ', 'በግ', 'በራዘዎ',
'አፍሪ', 'በአጠቃለሁ', 'ኤተናወቻ', 'ወሰጥ', 'አልጋ', 'ለሥነት', 'ወረት', 'የመቆየት', 'ወድድር', 'ኋው፡
', 'ደብብ', 'አፍሪ', 'ለበደብብ', 'ኋክ']
```

Size of valid tokens are: 80769388 Words(tokens)

```

In [7]: start_time = time.time()
tokens1 = nltk.word_tokenize(text)
end_time = time.time()
print(f"Time taken for tokenization: {(end_time - start_time)/60} minutes")

```

Time taken for tokenization: 12.268705590565999 minutes

Here, the entire file's content is read into one string, and then NLTK's word_tokenize function is used to tokenize the string into words. This method is straightforward and leverages the robust tokenization capabilities of the NLTK library. Now we have list of tokens in a single list variable `tokens1` but it is not `Efficient` so we use tokenization function `tokenizer` where built from scratch , because it tooks `12.22` minutes, almost 6 times customized one which is `2.24` minutes.

```

In [8]: print('Token size ',len(tokens))

```

Token size 80769388

```

In [9]: print('\n\n Sample  from the tokens: first 500 words are :\n')
print(tokens[:500])

```

Sample from the tokens: first 500 words are :

Save the tokens for later use and change the address to none temporary directory or download and manauly upload to fixed directory

The Kaggle working directory is temporary, and it deletes files upon restarting. To ensure the preservation of tokens for later use, one option is to save them and change the directory to a non-temporary location. Alternatively, the tokens can be downloaded and manually

uploaded to a fixed directory. In my case, I will manually upload the saved tokens to '/kaggle/input/saved-amharic-tokens'.

```
In [10]: # The Kaggle working directory is temporary, and it deletes files upon restarting.  
# To ensure the preservation of tokens for later use, one option is to save them as  
with open('/kaggle/working/tokens.pkl', 'wb') as file:  
    pickle.dump(tokens, file)
```

Load saved tokens we have to load the token from different directory

```
In [11]: # ## Load saved tokens we have to load the token from different directory
with open('/kaggle/input/saved-tokens/tokens.pkl', 'rb') as file:
    tokens = pickle.load(file)
```

Clean the corpus from amharic stopwords, English symbols, Numbers, amharic punctuations

After tokenization, the next step is to clean the corpus by removing stopwords, symbols, and punctuation from both the Amharic and English languages.

```
In [13]: # A sample List of custom stopwords in Amharic
import string

# Define a List of custom stopwords in Amharic
amharic_stopwords = []

# Define Amharic punctuations
amharic_punctuations = [':', '#', '%', '፣', '፤', '፤፡', '፤፤', '፤፤፡']

# English symbols that we want to remove (this includes letters, digits, and punctuation)
english_symbols = list(string.ascii_letters + string.digits + string.punctuation)

# Combined List of characters to remove
elements_to_remove = amharic_stopwords + amharic_punctuations + english_symbols
start_time = time.time()
cleaned_tokens = [token for token in tokens if not any(ch in elements_to_remove for ch in token)]
print('\nTokens removed \n', elements_to_remove)

end_time = time.time()
print(f"\n\nTime taken for cleaning : {((end_time - start_time)/60} minutes \n")
```

Tokens removed

[':', '#', '\$', '%', '&', '^', '(', ')', '*', '+', ',', '.', '/', ':', ';', '<', '=', '>', '?', '@', '[', '\\\\', ']', '^', - , '}', '|', '{', '~']

Time taken for cleaning : 9.963485936323801 minutes

First we define the unwanted elements, including Amharic stopwords, Amharic punctuation, and English symbols (comprising letters, digits, and punctuation). The then iterates through

the list of tokens and excludes any token that contains these elements or is composed solely of digits. The resulting cleaned_tokens list contains the text data cleansed of irrelevant characters and symbols, ready for further text analysis or processing tasks.

It then cleans the tokenized data by removing any tokens that contain these defined characters or are purely numeric. The result is a list of cleaned tokens free from unwanted characters.

```
In [14]: print('\nToken size after cleaned: ', len(cleaned_tokens))
print('\nThe first 50 words from the tokens \n',cleaned_tokens[:150])

print('\nToken size before cleaned: ',len(tokens))
print('\nThe first 50 words from the original tokens \n',tokens[:150])

print('\nLength difference/Number of removed tokens: ',len(tokens)-len(cleaned_tokens))
```

Token size after cleaned: 71813265

The first 50 words from the tokens

['ምን', 'መልታችሁ', 'አንበረዥን', 'አትዋኝያ', 'በተደረገማ', 'ጥርወ', 'ደርሱት', 'ልተታደሙው', 'የፈልግቶ
ወ', 'የከለም', 'የእጋር', 'ክብ', 'ዋ', 'ዘዘ', 'በደቦብ', 'አፍሪካ', 'በፌቅ', 'አየዋቻ', 'አንድታ', 'በረሬ,
ልክ', 'በኢመቱ', 'በለስ', 'ቀፍቃና', 'ለለ', 'ዋ', 'ልተታደም', 'ሁለት', 'ልቃቃና', 'ወደ', 'ደብብ', 'አፍ
ሪካ', 'ቢዋ', 'በረከብ', 'አፍሪካ', 'አብር', 'የመኖር', 'ወደድር', 'በደቦብ', 'አፍሪካ', 'ከተላያየ', 'የእና
ካ', 'አንድታ', 'የተውጣጠ', 'የህል', 'ተሳታፊዎች', 'የተከራለበት', 'ይህ', 'ግለሰቦች', 'በረታና', 'ሁኔታ',
'መሰጥ', 'በሚለች', 'በቃቃቃውን', 'የሚደረሰበትናበበት', 'መሆኑን', 'የሚገጥሚቃቃውን', 'የተለያየ', 'ኤተናዋች',
'በተካግባኝተና', 'በተበብ', 'ከለመቻ', 'ዘር', 'ተሰማցወች', 'ቃጊርቻን', 'በበልጠት', 'መጽታት', 'ወዘተ', 'በየ
ዘው', 'ከሚደረገው', 'ቁሳ', 'ተጨረወ', 'ቀናት', 'የህል', 'በወደድሩ', 'መቆሻት', 'የቃሉ', 'ሁለት', 'ተመዳ
ገሪዋች', 'አያንዲንቃቃው', 'ሽህ', 'ዶላር', 'አንድማሽለመ-ም', 'አናገር', 'በዘንድርው', 'ወደድር', 'አገራቻን',
'የኋ', 'አኋ', 'የኋ', 'የተባለ', 'ሁለት', 'መስጥቻን', 'በተሰራፍም', 'የኋ', 'ቀድሞ', 'በለው', 'የቀኑው',
'ሰላም', 'በሆነ', 'የኋዎ', 'በቅርቡ', 'ከወደድር', 'ወው', 'የእንዳቱ', 'በቃና', 'ተሰራ', 'ወደ', 'አገኘና
ት', 'ለማሻጋገር', 'የህዝብ', 'የድጋፍ', 'ድም', 'መሳና', 'መሆኑን', 'የተገኘበው', 'መቆጃ', 'ነው', 'አገኘና
ሁ', 'የሚል', 'አገራዊ', 'የኋ', 'የስተላፈልጉ', 'የኋ', 'የኋ', 'ከወደድሩ', 'ከመሰብበት', 'የእንደን', 'በም',
'በለስናነት', 'የሚከበረት', 'ከፍተኛ', 'አገራ', 'አንድታብና', 'የሚል', 'ከፍተኛ', 'በቃና', 'አንድራቢት', 'የ
ንቱ', 'የይሆን', 'በቃና', 'የወደድሩ', 'የዋሁት', 'አትዋኝያ', 'ይስና', 'ይሁን', 'የሚል', 'አንድሁ', 'ከን
ቱ', 'በቃና']

Token size before cleaned: 80769388

The first 50 words from the original tokens

['ምን', 'መለቻሁ', 'አንበረን', 'አቶኋያዥ', 'በተደራማዊ', 'ጥርወ', 'ደርሃት', 'ልተታደሙዋ', 'የፊልግለቻ', 'የከለም', 'የእጊር', 'ክስ', 'ወ', 'ለ19ኛ', 'ንዘ', 'በደብብ', 'አፍሪካ', 'ሰጠዋ', 'በኩቅ', 'አየች', 'አንድቻ', 'ባረራ', 'ልክ', 'በአመቱ', 'በለስ', 'ቀናትና', 'ለላ', 'ወ', 'ልተታደም', 'ሁለት', 'ልቻቸ', 'መዳ', 'ደቦብ', 'አፍሪካ', 'ለከት::ይናው', 'ቢግ', 'በረዘዘሪለ', 'አፍሪካ', 'ከነጋ', 'የመሞር', 'መድናጋ', 'በደብብ', 'አፍሪካ', 'ቴክኖሎጂ::', 'ከተለያየ', '14', 'የአፍሪካ', 'አገራት', 'የተወጠጣ', '26', 'የህል', 'ተሳታፊዎች', 'የተከፈለበት', 'ይህ', 'ወድድርጋ', 'ግለሰቦች', 'በረታኝ', 'ሁኔታ', 'ወስት', 'በማለፍ', 'በቻታቻውን', 'የሚረስመሰከሩበት', 'መሆኑ', 'አምተና::', 'የሚገጥሟችውን', 'የተለያየ', 'ፈተናምች', 'በተካሳት', 'በተበብ', 'ማለፍ::', 'ከለመቻች', 'የጋር', 'ተስማምዎ', 'መዘላቹ', 'ቁጥርቻች', 'በተገኘው', 'በወጪዎች', 'መቆጥት', 'መዘተ', 'የበረከው', 'ከሚደረገው', 'ቁና', 'ተርጉሙ', 'ለ91', 'ቀናት', 'የህል', 'በወድድሩ', 'መቆጥት', 'የጽኢሉ', 'ሁለት', 'ተወካከለዋች', 'አንድገትናው', '200', 'ይህ', 'ዶክ', 'አንድሚሰመው', 'ልኝነት', 'ነበር::', 'በዘንድርው', 'ወድድር', 'አገራቻ', 'ክን', 'የኅ', 'የተባለ', 'ሁለት', 'ወጣዋች', 'በታሰቢው', 'የኅ', 'ቀደም', 'በለው', 'የቀኑው', 'በለው', 'የቀኑው', 'ለለዋ', 'ለሁን', 'የነዋም', 'ለቀርቡ', 'ከወደድር', 'ወመ', 'ሆናቸውን', 'የአገራቱ', 'በቻች', 'ተስፋ', 'መዳ', 'አገናዘነት', 'ለማሽንጂ', 'የህዝብ', 'የድጋፍ', 'ይም', 'መናና', 'መሆኑ', 'የተገኘበው', 'መደኝ', 'ነው', 'አንጋዥ', '835', 'የሚል', 'አገራዊ', 'ጥር', 'የሰተላለፈኝ', 'የኅ', 'የኅ', 'ከወደድሩ', 'ከመሰበበት', 'በረቱ::መደኝ', 'የአገኝ', 'ለም', 'በአገናዘነት', 'የሚሰጠቸት', 'ከፍተቻ', 'ገንዘብ', 'አገኝ']

Length difference/Number of removed tokens: 8956123

Save the cleaned tokens for later use

```
In [15]: # Save the tokens for later use
with open('/kaggle/working/tokens_cleaned.pkl', 'wb') as file:
    pickle.dump(cleaned_tokens, file)
```

```
In [16]: # # Load saved clean_Tokens, we have to Load the token from different directory  
  
with open('/kaggle/input/saved-tokens/tokens_cleaned(1).pkl', 'rb') as file:  
    tokens = pickle.load(file)  
  
# with open('/kaggle/working/tokens_cleaned.pkl', 'rb') as file:  
#     tokens = pickle.load(file)
```

N-gram language mode

1.1 Create n-grams for n=1, 2, 3, 4.

- Top

```
In [19]: # Extract n-grams n-gram generators and change to list
start_time = time.time()
unigrams = ngrams(tokens, 1)
bigrams = ngrams(tokens, 2)
trigrams = ngrams(tokens, 3)
fourgrams = ngrams(tokens, 4)
end_time = time.time()
print(f"\n\nTime taken for n-gram generation : {(end_time - start_time)/60} minutes

# we can change it to list but it is not memory efficient
unigramslist = list(ngrams(tokens, 1))
bigramslist = list(ngrams(tokens, 2))
trigramslist = list(ngrams(tokens, 3))
fourgramslist = list(ngrams(tokens, 4))
```

Time taken for n-gram generation : 2.551078796386718e-06 minutes

```
In [20]: print(fourgramslist[:10])
```

Print samples first 50 n-grams for all 4 n-grams

```
In [21]: print('\n\nFirst 50 Unigrams\n')
print(unigramslist[:50])
print('\n')
print('\nFirst 50 Bigrams\n')
print(bigramslist[:50])
print('\n')
print('\nFirst 50 Trigrams\n')
print(trigramslist[:50])
```

```
print('\n')
print('\nFirst 50 fourgrams\n')
print(fourgramslist[:55])
```

First 50 Unigrams

First 50 Bigrams

[('ምን', 'መስራቶሁ'), ('መስራቶሁ', 'አንበያን'), ('አንበያን', 'አትናዕም'), ('አትናዕም', 'በተዳደሪዎች'), ('በተዳደሪዎች', 'ጥጋው'), ('ጥጋው', 'ደርሱት'), ('ደርሱት', 'ልተታደሙዋ'), ('ልተታደሙዋ', 'የልጅለቶዎ'), ('የልጅለቶዎ', 'የእለም'), ('የእለም', 'የእግዚ'), ('የእግዚ', 'ክቡ'), ('ክቡ', 'ዋ'), ('ዋ', 'ዘዴ'), ('ዘዴ', 'በደብብ'), ('በደብብ', 'አፍሪካ'), ('አፍሪካ', 'በኩቅ'), ('በኩቅ', 'አያቶች'), ('አያቶች', 'አንድጋድ'), ('አንድጋድ', 'ባረሪ'), ('ባረሪ', 'ልክ'), ('ልክ', 'በአመቱ'), ('በአመቱ', 'በላሉ'), ('በላሉ', 'ቀናቴና'), ('ቀናቴና', 'ለለ'), ('ለለ', 'ዋ'), ('ዋ', 'ልተታደሙዋ'), ('ልተታደሙዋ', 'ሁለት'), ('ሁለት', 'ልፋይን'), ('ልፋይን', 'ወደ'), ('ወደ', 'ደብብ'), ('ደብብ', 'አፍሪካ'), ('አፍሪካ', 'በግ'), ('በግ', 'ብርሃን'), ('ብርሃን', 'አፍሪካ'), ('አፍሪካ', 'አብር'), ('አብር', 'የመኖር'), ('የመኖር', 'ውድድር'), ('ውድድር', 'በደብብ'), ('በደብብ', 'አፍሪካ'), ('አፍሪካ', 'ከተለያቶ'), ('ከተለያቶ', 'የእናዕም'), ('የእናዕም', 'የእናዕም'), ('የእናዕም', 'አገልታ'), ('አገልታ', 'የተውጣሁ'), ('የተውጣሁ', 'የህል'), ('የህል', 'ተሳታፊዎች'), ('ተሳታፊዎች', 'የተከፈልሁት'), ('የተከፈልሁት', 'የህሉ'), ('የህሉ', 'በፈታኝ'), ('በፈታኝ', 'ሁኔታ')]]

First 50 Trigrams

First 50 fourgrams

```
In [22]: # sample prints
# print samples for all 4 n-grams
print('\nUnigrams')
print(next(unigrams))
print(next(unigrams))
print(next(unigrams))
print(next(unigrams))

print('\nBigrams')
print(next(bigrams))
print(next(bigrams))
print(next(bigrams))
print(next(bigrams))

print('\nTrigrams')
print(next(trigrams))
print(next(trigrams))
print(next(trigrams))
print(next(trigrams))
print(next(trigrams))

print('\nFourgrams')
print(next(fourgrams))
print(next(fourgrams))
print(next(fourgrams))
print(next(fourgrams))
```

Unigrams

('ምኑ' ,)
('መስተካከል' ,)
('እንበብያን' ,)
('እትዮጵያ' ,)

Bigrams

('ምኑን', 'መስላለችሁ')
 ('መስላለችሁ', 'አንበበያን')
 ('አንበበያን', 'አትዮጵያ')
 ('አትዮጵያ', 'በተደረገመሳቅ')

Trigrams

(‘ምን’, ‘መስራቶሁ’, ‘እንዲበደን’) (‘መስራቶሁ’, ‘እንዲበደን’, ‘እተናቋያ’) (‘እንዲበደን’, ‘እተናቋያ’, ‘በተደረጋጋሚ’) (‘እተናቋያ’, ‘በተደረጋጋሚ’, ‘ጥረዋ’)

Fourgrams

- (‘ምን’, ‘መለያቸው’, ‘አንበሳን’, ‘አቶዕዳደሪ’)
- (‘መለያቸው’, ‘አንበሳን’, ‘አቶዕዳደሪ’, ‘በተደረገማ’)
- (‘አንበሳን’, ‘አቶዕዳደሪ’, ‘በተደረገማ’, ‘ጥራው’)
- (‘አቶዕዳደሪ’, ‘በተደረገማ’, ‘ጥራው’, ‘ደርሃት’)

1.2. Calculate probabilities of n-grams and find the top 10 most likely n-grams for all n.

- Top

1.2.1. probabilities of 2-grams and top 10 most likely 2-grams

Calculate probabilities

```
In [23]: # Get frequency of bigrams
bigram_freq = Counter(bigrams)
trigram_freq = Counter(trigrams)
fourgram_freq = Counter(fourgrams)

# # Calculate probabilities
total_bigrams = sum(bigram_freq.values())
bigram_probabilities = {gram: freq/total_bigrams for gram, freq in bigram_freq.items()}
```

Get top 10 bigrams

```
In [24]: top_10_bigrams = sorted(bigram_probabilities.items(), key=lambda item: item[1], reverse=True)
print("{:<40} {:<20}".format("\nTop 10 fourgrams ", " Probability"))

# Print each element in the desired format
for ngram, probability in top_10_bigrams:
    ngram_str = ' '.join(ngram)
    print("{:<40} {:<20}".format(ngram_str, probability))
```

Top 10 fourgrams	Probability
א א	0.0018562964337997358
ה ה	0.0007995556147818973
ו ו	0.0007798519096186599
ב ב	0.00041674077161042754
ל ל	0.00038229632461454256
ה ה	0.00037214817571467967
ו ו	0.00037207410163511863
ל ל	0.00036607410119067415
ה ה	0.00035088891488066037
ו ו	0.0003265926167846383

1.2.2. probabilities of trigrams and top 10 most likely trigrams

Calculate probabilities

```
In [25]: total_trigrams = sum(trigram_freq.values())
trigram_probabilities = {gram: freq/total_trigrams for gram, freq in trigram_freq.items()}
```

Get top 10 trigrams

```
In [26]: top_10_trigrams = sorted(trigram_probabilities.items(), key=lambda item: item[1], reverse=True)
print("{:<40} {:<20}".format("\nTop 10 fourgrams ", " Probability"))

# Print each element in the desired format
for ngram, probability in top_10_trigrams:
    ngram_str = ' '.join(ngram)
    print("{:<40} {:<20}".format(ngram_str, probability))
```

Top 10 fourgrams	Probability
ՓԴ Գ Թ	0.000776148263133076
Հ Հ Հ	0.00034837042198080325
Գ Թ ՔՊԸ	0.00011970372143758836
ԺՄԻՐՇ ԹՄՀԸ ԱՌ	0.00010162964468587329
ՀՈՊ ԻՒՄ ՀՈՒՅԸԸԸ	8.792593895199096e-05
ՔՀԶՆ ՀՈՊ ԻՒՄ	8.214816031824597e-05
ՈՊԼՊ ՀՓԲ ՋՀԸ	6.70370469684514e-05
ԻՆԻ ՄՁ ՂԱ	6.548149118244314e-05
ՊԸ ՓԴ Գ	6.540741709739512e-05
ՔԻՒՇ ՓԴ Գ	6.496297258710705e-05

1.2.3. probabilities of 4-grams and top 10 most likely 4-grams

Calculate probabilities

```
In [27]: total_fourgrams = sum(fourgram_freq.values())
fourgrams_probabilities = {gram: freq/total_fourgrams for gram, freq in fourgram_fr
```

Get top 10 fourgrams

```
In [28]: top_10_fourgrams = sorted(fourgrams_probabilities.items(), key=lambda item: item[1])
print("{:<40} {:<20}".format("\nTop 10 fourgrams ", " Probability"))

# Print each element in the desired format
for ngram, probability in top_10_fourgrams:
    ngram_str = ' '.join(ngram)
    print("{:<40} {:<20}".format(ngram_str, probability))
```

Top 10 fourgrams	Probability
ՊԸ ՓԴ Գ Թ	6.503705148971514e-05
ՔԻՒՇ ՓԴ Գ Թ	6.48148292181102e-05
ՄՊՇՈՒՇ ՓԴ Գ Թ	6.25926065020607e-05
ԴՓՄՇՆ ՓԴ Գ Թ	5.31112291358287e-05
ՊՂՈՒՇ ՓԴ Գ Թ	5.222223382716307e-05
ՀՄՈԼԻՒՇ ՈՒՇ ԲՊԱՐՅՆԹՈՎԸ ՄՊՄՒՇ	5.214815973662809e-05
ՔՀԺԺՄՆ ՀՄՈԼԻՒՇ ՈՒՇ ԲՊԱՐՅՆԹՈՎԸ	5.029630747325351e-05
ՑԺԿ ՔՀԺԺՄՆ ՀՄՈԼԻՒՇ ՈՒՇ	5.0148159292183544e-05
Գ Գ Գ Թ ՔՊԸ	4.977778883950863e-05
ՈՆ ՓԴ Գ Թ	4.851852930041392e-05

1.3. What is the probability of the sentence. "ՀԵՐՔԸ ԺՇԽՎ ՄԴՀ ՀԴԴ".

- Top

```
In [29]: def generate_ngrams(tokens, n):
    return list(ngrams(tokens, n))

#n-gram_probabilities is a dictionary containing n-gram probabilities
def probability_of_sentence(sentence,n):
    sentence_tokens = sentence.split()
    ngrams_sentence = generate_ngrams(sentence_tokens, n)
    sentence_probability = 1.0

    if n==2:
        for bg in ngrams_sentence:
```

```

        if bg in bigram_probabilities:
            sentence_probability *= bigram_probabilities[bg]
        else:
            sentence_probability *= 0.0 # Assuming we assign zero probability
    return sentence_probability
elif n==3:
    for bg in ngrams_sentence:
        if bg in trigram_probabilities:
            sentence_probability *= trigram_probabilities[bg]
        else:
            sentence_probability *= 0.0 # Assuming we assign zero probability
    return sentence_probability
elif n==4:
    for bg in ngrams_sentence:
        if bg in fourgrams_probabilities:
            sentence_probability *= fourgrams_probabilities[bg]
        else:
            sentence_probability *= 0.0 # Assuming we assign zero probability
    return sentence_probability
else:
    return 'Please input n 2, 3, or 4'

sentence1 = "አትወቻ ቤትዎች የገዢ"

print(f"\n\n\nThe probability of the sentence(አትወቻ ቤትዎች የገዢ የገዢ) using bigram is:")
print(f"\nThe probability of the sentence(አትወቻ ቤትዎች የገዢ የገዢ) using trigram is:")
print(f"\nThe probability of the sentence(አትወቻ ቤትዎች የገዢ የገዢ) using fourgram is:")

```

The probability of the sentence(አትወቻ ቤትዎች የገዢ የገዢ) using bigram is: 4.389575735
4063326e-20

The probability of the sentence(አትወቻ ቤትዎች የገዢ የገዢ) using trigram is: 5.48697007
5700204e-15

The probability of the sentence(አትወቻ ቤትዎች የገዢ የገዢ) using fourgram is: 7.4074090
53498309e-08

Bigram Model (n=2) : The probability is incredibly low at approximately 4.39e-20. This suggests that, according to the bigram model, the chance of this exact sequence of words appearing together is extremely rare within the context of the text corpus used for the model.

Trigram Model (n=3) : The probability increases significantly to around 5.49e-15 when using a trigram model. This implies that considering a sequence of three words at a time provides a higher likelihood for this sentence to be observed, reflecting a more common occurrence in the source material.

Fourgram Model (n=4) : The probability surges to approximately 7.41e-08 when evaluated with a fourgram model. This indicates that the sequence of four words is much more common, and the sentence is much more likely to occur in the dataset that was used to train the model.

The increasing probabilities of the Amharic sentence "አትወቻ ቤትዎች የገዢ የገዢ" with larger n-grams from bigrams to fourgrams demonstrate the importance of context in language

modeling. The low probability with the bigram model indicates a rare occurrence of word pairs, while the higher probabilities with trigram and fourgram models suggest that longer word sequences from the sentence are more common in the dataset. This trend highlights the need for extensive and context-rich datasets to train accurate language models, especially for languages with complex structures like Amharic.

a) using Bigrams

```
In [30]: sentence1 = "አትኋኑ ታሪዎች ማስታወሻ ነት"
sentence2 = "በፈቻቸው ሁኔታ ወሰን ማለፍ"
sentence3 = 'ተወካይ ጥምር በት'
sentence4 = 'ቦግለም አቀፍ ይረዳ'
sentence5 = 'በአዲስ አበባ'

print(f"\n\n\nThe probability of the sentence(አትኋኑ ታሪዎች ማስታወሻ ነት) is: \
{probability_of_sentence(sentence1, 2)}")
print(f"\nThe probability of the sentence(በፈቻቸው ሁኔታ ወሰን ማለፍ) is: \
{probability_of_sentence(sentence2, 2)}")
print(f"\nThe probability of the sentence(ተወካይ ጥምር በት) is: \
{probability_of_sentence(sentence3, 2)}")
print(f"\nThe probability of the sentence(ቦግለም አቀፍ ይረዳ) is: \
{probability_of_sentence(sentence4, 2)}")
print(f"\nThe probability of the sentence(በአዲስ አበባ) is: \
{probability_of_sentence(sentence5, 2)})")
```

The probability of the sentence(አትኋኑ ታሪዎች ማስታወሻ ነት) is: 4.3895757354063326e-20

The probability of the sentence(በፈቻቸው ሁኔታ ወሰን ማለፍ) is: 1.0239904729428439e-17

The probability of the sentence(ተወካይ ጥምር በት) is: 4.169986351657763e-08

The probability of the sentence(ቦግለም አቀፍ ይረዳ) is: 1.7296891451391234e-08

The probability of the sentence(በአዲስ አበባ) is: 0.00036607410119067415

b) Using trigrams

```
In [31]: sentence1 = "አትኋኑ ታሪዎች ማስታወሻ ነት"
sentence2 = "በፈቻቸው ሁኔታ ወሰን ማለፍ"
sentence3 = 'ተወካይ ጥምር በት'
sentence4 = 'ቦግለም አቀፍ ይረዳ'

print(f"\n\nThe probability of the sentence(አትኋኑ ታሪዎች ማስታወሻ ነት) is: \
{probability_of_sentence(sentence1, 3)}")
print(f"\nThe probability of the sentence(በፈቻቸው ሁኔታ ወሰን ማለፍ) is: \
{probability_of_sentence(sentence2, 3)}")

print(f"\nThe probability of the sentence(ተወካይ ጥምር በት) is: \
{probability_of_sentence(sentence3, 3)}")
print(f"\nThe probability of the sentence(ቦግለም አቀፍ ይረዳ) is: \
{probability_of_sentence(sentence4, 3)})")
```

The probability of the sentence(አትኩክያ ቤርሃዊ ማንኛውም የተመለከተውን ነጥቦች) is: 5.486970075700204e-15

The probability of the sentence(በፈቻቸው ሁኔታ ወሰን መግለጫ) is: 0.0

The probability of the sentence(ተወካይነት ምክር በት) is: 0.00010162964468587329

The probability of the sentence(በዓለም አቀፍ ይረዳል) is: 6.70370469684514e-05

c) Using fourgrams

In [32]:

```
sentence1 = "አትኩክያ ቤርሃዊ ማንኛውም የተመለከተውን ነጥቦች"
sentence2 = "ገለሰቦች በፈቻቸው ሁኔታ ወሰን መግለጫ በቃቻቸውን የሚያስተካክሱት መሆኑን"
sentence3 = "በፈቻቸው ሁኔታ ወሰን መግለጫ"

print(f"\n\nThe probability of the sentence(አትኩክያ ቤርሃዊ ማንኛውም የተመለከተውን ነጥቦች) is: \
{probability_of_sentence(sentence1, 4)}")
print(f"\n\nThe probability of the sentence(ገለሰቦች በፈቻቸው ሁኔታ ወሰን መግለጫ በቃቻቸውን የሚያስተካክሱት መሆኑን) is: \
{probability_of_sentence(sentence2, 4)}")

print(f"\n\nThe probability of the sentence(በፈቻቸው ሁኔታ ወሰን መግለጫ) is: \
{probability_of_sentence(sentence3, 4)}")
```

The probability of the sentence(አትኩክያ ቤርሃዊ ማንኛውም የተመለከተውን ነጥቦች) is: 7.407409053498309e-08

The probability of the sentence(ገለሰቦች በፈቻቸው ሁኔታ ወሰን መግለጫ በቃቻቸውን የሚያስተካክሱት መሆኑን) is: 2.2301374979696553e-36

The probability of the sentence(በፈቻቸው ሁኔታ ወሰን መግለጫ) is: 7.407409053498309e-08

1.4. Generate random sentences using n-grams; explain what happens as n-increases based on your output

- [Top](#)

Define random sentence generator function

In [33]:

```
def sentence_generator(gram_freq,n, token, num_words=8):
    if n==2:
        bigram_sentence = [token]
        for _ in range(num_words - 1):
            next_tokens = [pair[1] for pair in gram_freq if pair[0] == token]
            if not next_tokens:
                break
            token = random.choice(next_tokens)
            bigram_sentence.append(token)
        return ' '.join(bigram_sentence)
    elif n==3:
        trigram_sentence = [token]
        for _ in range(num_words - 1):
            next_tokens = [pair[1] for pair in gram_freq if pair[0] == token]
            if not next_tokens:
                break
            token = random.choice(next_tokens)
            trigram_sentence.append(token)
        return ' '.join(trigram_sentence)
    elif n==4:
```

```
fourgram_sentence = [token]
for _ in range(num_words - 1):
    next_tokens = [pair[1] for pair in gram_freq if pair[0] == token]
    if not next_tokens:
        break
    token = random.choice(next_tokens)
    fourgram_sentence.append(token)
return ' '.join(fourgram_sentence)
else:
    return 'Please input one of the following n values n=2,n=3, n=4'
```

a) Using bigrams

```

In [34]: import random
# def generate_sentence(bigram_freq, token, num_words=8):
#     sentence = [token]
#     for _ in range(num_words - 1):
#         next_tokens = [pair[1] for pair in bigram_freq if pair[0] == token]
#         if not next_tokens:
#             break
#         token = random.choice(next_tokens)
#         sentence.append(token)
#     return ' '.join(sentence)

# Generate a sentence starting with a given word
starting_word1 = 'አትዋኝ'
starting_word2 = 'በተደረጋሚ'
starting_word3 = 'ከተለያየ'
starting_word4 = 'የእለም'
starting_word5 = 'በደብብ'
starting_word6 = 'የእናርሻ'
starting_word7 = 'በፋይ'

print('አትዋኝ as a starting word = ', sentence_generator(bigram_freq,2, starting_word1))
print('በተደረጋሚ as a starting word = ', sentence_generator(bigram_freq,2, starting_word2))
print('ከተለያየ as a starting word = ', sentence_generator(bigram_freq,2, starting_word3))
print('የእለም as a starting word = ', sentence_generator(bigram_freq,2, starting_word4))
print('በደብብ as a starting word = ', sentence_generator(bigram_freq,2, starting_word5))
print('የእናርሻ as a starting word = ', sentence_generator(bigram_freq,2, starting_word6))
print('በፋይ as a starting word = ', sentence_generator(bigram_freq,2, starting_word7))

አትዋኝ as a starting word = አትዋኝ ወጪው ጥለያቻች የቃኝ ደንበር መመለስ ያልፈገን እኔ
በተደረጋሚ as a starting word = በተደረጋሚ የሚችቁዸውንታና መድቦች ጥነዚው በቀረበው የሚሆን አጥጋቢን ፕሏ
ከተለያየ as a starting word = ከተለያየ ምግባር ወሰጥ አድራሻት አነዱች ድግሞ ምናለሁት ይህንን
የእለም as a starting word = የእለም ቁርቡችን ባካላ የአትዋኝውያን በሁሉ ልማድ ከተጨማሪው የእክምሮ
በደብብ as a starting word = በደብብ አጭርቃው ማኅበዋል ወይም የእናርሻ ተረዘዋንት አንዳንድ የበረት
የእናርሻ as a starting word = የእናርሻ ቁጥሮች መመዘኛ ተደርሱ አስተዳደሩን አለሁ ለሰነድበቻው
በፋይ as a starting word = በፋይ ስምተኞች በሌሎች ማስወጣ አገማዎችን ወጪ ቁጥሮች ወጪ ይ

```

b) Using trigrams

```
In [35]: # Generate a sentence starting with a given word
starting_word1 = 'አትዋኑ'
starting_word2 = 'በተዳደሪያዎች'
starting_word3 = 'ከተለያየ'
starting_word4 = 'የካለም'
starting_word5 = 'በደህንብ'
starting_word6 = 'የእናፍቅ'
```

```
starting_word7 = 'ԱՌՎ'
```

```
print('ՀԴՔՋ as a starting word = ', sentence_generator(trigram_freq, 3, starting_w
print('ՈՒԾՉՉՄՆ as a starting word = ', sentence_generator(trigram_freq, 3, starting_
print('ԻՒԱՐՔ as a starting word = ', sentence_generator(trigram_freq, 3, starting_wor
print('ՔԱԼՊՈ as a starting word = ', sentence_generator(trigram_freq, 3, starting_wor
print('ՈԶՈՒՆ as a starting word = ', sentence_generator(trigram_freq, 3, starting_wo
print('ՔՀԳՃԻ as a starting word = ', sentence_generator(trigram_freq, 3, starting_wor
print('ՈՒՎ as a starting word = ', sentence_generator(trigram_freq, 3, starting_word)
```

ՀԴՔՋ as a starting word = ՀԴՔՋ իսպոյդիշն ԱՊՐԵ ԱԼԱ ԺԻՆ ԳԵ ԶՊՊՈ ՈԽԱԺՄՈՅՔԴՈ
ՈՒԾՉՉՄՆ as a starting word = ՈՒԾՉՉՄՆ ԹՀՀ ԳՈՒ ՔՄԱԲԴԻՆ ՈՄԱԼՈՒ ՏՊԱԿՈԸ ՄՊՄԻՆ ՈՂՄ-ՔԳԸ
ԻՒԱՐՔ as a starting word = ԻՒԱՐՔ ՀԻՂՈՄՔ ՄՈՒ ԻՆՉԻ ՈՐԴ ԱՅ ՄՊԻՆՔ ՔՄԱԾՄԸ
ՔԱԼՊՈ as a starting word = ՔԱԼՊՈ ՄԸՆԴ ԲՍԱ ԴՈՒՆ ԱՆՁԳԵՇ ՈՄՑԺԿ ՄՈՒՐ ՄՊՈԴ
ՈԶՈՒՆ as a starting word = ՈԶՈՒՆ ԱՀՈ ԿՆ ԻՄԱԼՈՒ ՄՈՒՐ ՄԿՆ ՓՀՆ ՈՂՄ
ՔՀԳՃԻ as a starting word = ՔՀԳՃԻ ՀԿԸ ՔԱԺՄՈՎ ԴԳՐԵՇ ՀՄ-ՔԸ ՄՈԽՈՒ ՍԱՄ ՈՂՄԱՊՐԵԴ
ՈՒՎ as a starting word = ՈՒՎ ՈՄՔ ՄՔ ՀԵՊ ՈՄՄԻՆ ԱՆՁՆ ՔՀԸ ՄԺԸԾԵՇ

c) Using fourgrams

In [36]:

```
starting_word1 = 'ՀԴՔՋ'
starting_word2 = 'ՈՒԾՉՉՄՆ'
starting_word3 = 'ԻՒԱՐՔ'
starting_word4 = 'ՔԱԼՊՈ'
starting_word5 = 'ՈԶՈՒՆ'
starting_word6 = 'ՔՀԳՃԻ'
starting_word7 = 'ՈՒՎ'
```

```
print('ՀԴՔՋ as a starting word = ', sentence_generator(fourgram_freq, 4, starting_w
print('ՈՒԾՉՉՄՆ as a starting word = ', sentence_generator(fourgram_freq, 4, starting_
print('ԻՒԱՐՔ as a starting word = ', sentence_generator(fourgram_freq, 4, starting_w
print('ՔԱԼՊՈ as a starting word = ', sentence_generator(fourgram_freq, 4, starting_w
print('ՈԶՈՒՆ as a starting word = ', sentence_generator(fourgram_freq, 4, starting_w
print('ՔՀԳՃԻ as a starting word = ', sentence_generator(fourgram_freq, 4, starting_w
print('ՈՒՎ as a starting word = ', sentence_generator(fourgram_freq, 4, starting_word)
```

ՀԴՔՋ as a starting word = ՀԴՔՋ ՓՄՈՆ ՈՐԵԴՔՆԵԿ ՔՄԱՄԱԼԱՎ ԹՄԴԲՆՈՎ ՔԼՎ ԺԸՎ ՀԴՔՋԸ
ՈՒԾՉՉՄՆ as a starting word = ՈՒԾՉՉՄՆ ԴՄՈՒՆԵ ՀՔՄ ԻՔԻ ՄԶԱ ՈՒՄՊԵՄ ՓԸՎԸ ՈՓԸՎԸ
ԻՒԱՐՔ as a starting word = ԻՒԱՐՔ ՔԲԴԻ ԳԵ ԱԾ ԱՅ ՔԱՓ ՖԱ ՈՒՄԴԵՄ
ՔԱԼՊՈ as a starting word = ՔԱԼՊՈ ԺԱՂՎ ՀԴՔՋՔՄՈՎ ՄԴՅԱՊՄԵՄ ԿՆ ԱԼՄ ԱՄԸՄ ԱՄԸՄ-ՔՄ
ՈԶՈՒՆ as a starting word = ՈԶՈՒՆ ԱՀՈ ԱՐԴՊՀՀ ԸՍ ԱԿ ՓՈԱԹՔ ՄՈՒՐ ՔԱ
ՔՀԳՃԻ as a starting word = ՔՀԳՃԻ ՀԿԸ ՄՊՀՈՒ ՊՀՆ ԳԵՆ ԳԵ ՈՒԱԾ ՔՀԳՃԻ ՀԿԸ
ՈՒՎ as a starting word = ՈՒՎ ԱԼՄ ՀՈՎՃ ՄՈՒ ԿՄՄԵԴ ՄՈԱԿՄ ՀՈՒԹԵՇՆ ՀԴՊՂՆՎ

2. Evaluate these Language Models Using Intrinsic Evaluation Method.

Perplexity

In order to implement the perplexity formula, you'll need to know how to implement m-th order root of a variable.

$$PP(W) = \sqrt[m]{\prod_{i=1}^m \frac{1}{P(w_i|w_{i-1})}}$$

From calculus:

$$\sqrt[M]{\frac{1}{x}} = x^{-\frac{1}{M}}$$

Here is bellow a code that will help us with the formula.

- [Top](#)

```
In [37]: import math
from collections import Counter

def calculate_perplexity(ngrams, total_words):
    model_entropy = 0.0
    for ngram, count in ngrams.items():
        prob = count / total_words
        model_entropy += -math.log2(prob)

    model_entropy /= total_words
    return math.pow(2, model_entropy)

def evaluate_model(n, train_tokens, test_tokens):

    ngrams_train = Counter(tuple(train_tokens[i:i + n]) for i in range(len(train_to
total_words_train = len(train_tokens))

    ngrams_test = Counter(tuple(test_tokens[i:i + n]) for i in range(len(test_to
total_words_test = len(test_tokens)

    return calculate_perplexity(ngrams_test, total_words_test)
# Load the tokens and calculate frequencies
with open('/kaggle/input/saved-tokens/tokens_cleaned(1).pkl', 'rb') as file:
    tokens = pickle.load(file)
# Split the data into training and test sets (80% training, 20% test)
train_size = int(0.8 * len(tokens))
train_tokens = tokens[:train_size]
test_tokens = tokens[train_size:]

# Evaluating n-gram models for n = 1, 2, 3, 4
perplexities = []
for n in range(1, 5):
    perplexity = evaluate_model(n, train_tokens, test_tokens)
    perplexities.append(perplexity)

# Storing the results in a list of tuples
results = []

for n, perplexity in enumerate(perplexities):
    results.append((n + 1, perplexity))

# Printing the results
for n, perplexity in results:
    print(f"Perplexity of {n}-gram model is : {perplexity}")
```

Perplexity of 1-gram model is : 3.4376775176273546
Perplexity of 2-gram model is : 13141.287652919986
Perplexity of 3-gram model is : 651342.603025604
Perplexity of 4-gram model is : 1549724.6862948143

The perplexity values you obtained indicate how well our n-gram models generalize to unseen data. Perplexity is a measure of uncertainty or "surprise" associated with predicting the next word in a sequence. A lower perplexity indicates better performance.

Here's a brief interpretation of the results:

Unigram (n=1):

Perplexity: 3.44 The unigram model has relatively low perplexity, suggesting that it provides a good estimate of word probabilities based on individual words. This is expected, as unigrams consider each word in isolation.

Bigram (n=2):

Perplexity: 13141.29 The bigram model has a significantly higher perplexity compared to the unigram model. This indicates that predicting the next word based on the previous word alone might not be as accurate. It suggests that the model struggles with capturing dependencies between adjacent words.

Trigram (n=3):

Perplexity: 651342.60 The trigram model has an even higher perplexity, indicating that considering the previous two words for prediction does not lead to effective generalization. The model might be overfitting to the training data and failing to capture broader linguistic patterns.

Fourgram (n=4):

Perplexity: 1549724.69` The fourgram model exhibits the highest perplexity, suggesting that considering the previous three words for prediction may lead to overfitting or insufficient data for effective generalization. The model might struggle with unseen combinations of four words.

Description: The results highlight the trade-off between model complexity and generalization. While the unigram model performs well, higher-order n-gram models (bigram, trigram, fourgram) face challenges in capturing dependencies and suffer from increasing perplexity, indicating potential overfitting to the training data. This could be due to the limitations of n-gram models in handling the complexity of natural language, especially when dealing with longer dependencies between words. Consideration of more advanced language models or techniques may be warranted for improved performance on this task.

3. Evaluate these Language Models Using Extrinsic Evaluation Method

- [Top](#)

Create predictive model for extrinsic evaluation

```
In [38]: # Function for predictive text entry (for extrinsic evaluation)
def predict_next_word(current_sentence, bigram_freq):
    current_tokens = current_sentence.split()
    last_token = current_tokens[-1]
    next_words = {pair[1]: freq for pair, freq in bigram_freq.items() if pair[0] == last_token}
    next_word = max(next_words, key=next_words.get) if next_words else None
    return next_word
```

Extrinsic evaluation - predict next word

```
In [39]: # Extrinsic evaluation - predict next word
current_sentence = "ታደሮ ህንጻ"
next_word = predict_next_word(current_sentence, bigram_freq)
print('\n\nThe previous words are: ', current_sentence)
print(f'\n and the suggested word is: {next_word}')
```

The previous words are: ተደሮ ህንጻ

and the suggested word is: ወሳኑ

```
In [40]: current_sentence = "ታደሮ"
next_word = predict_next_word(current_sentence, bigram_freq)
print('\n\nThe previous words are: ', current_sentence)
print(f'\n and the suggested word is: {next_word}')
```

The previous words are: ተደሮ

and the suggested word is: የሸ

```
In [41]: current_sentence = "በፌዴራል ማና"
next_word = predict_next_word(current_sentence, bigram_freq)
print('\n\nThe previous words are: ', current_sentence)
print(f'\n and the suggested word is: {next_word}')
```

The previous words are: በፌዴራል ማና

and the suggested word is: ወሳኑ

```
In [42]: current_sentence = "ውሳኑ በማለፍ"
next_word = predict_next_word(current_sentence, bigram_freq)
print('\n\nThe previous words are: ', current_sentence)
print(f'\n and the suggested word is: {next_word}')
```

The previous words are: ወሳኑ በማለፍ

and the suggested word is: ወያ

Conclusion

The model utilizes bigram frequencies to predict the next word in a sentence, and the given test cases show that it can suggest contextually reasonable continuations. For instance, from the historical context indicated by "ታደሮ ህንጻ," the model aptly predicts "ውሳኑ," reflecting a coherent phrase structure in the language being tested. The other test cases align similarly with the expected linguistic patterns. Thus, based on these examples, we can infer that the

model is fairly effective in its predictive capability, indicating a good understanding of word pair probabilities within the corpus from which the bigram frequencies were derived.